

SysScale: Exploiting Multi-domain Dynamic Voltage and Frequency Scaling for Energy Efficient Mobile Processors

Jawad Haj-Yahya[§] Mohammed Alser[§] Jeremie Kim[§] A. Giray Yağlıkçı[§]
 Nandita Vijaykumar^{§*†} Efraim Rotem^{*} Onur Mutlu[§]

[§]ETH Zürich

^{*}Intel

[†]University of Toronto

There are three domains in a modern thermally-constrained mobile system-on-chip (SoC): compute, IO, and memory. We observe that a modern SoC typically allocates a fixed power budget, corresponding to worst-case performance demands, to the IO and memory domains even if they are underutilized. The resulting unfair allocation of the power budget across domains can cause two major issues: 1) the IO and memory domains can operate at a higher frequency and voltage than necessary, increasing power consumption and 2) the unused power budget of the IO and memory domains cannot be used to increase the throughput of the compute domain, hampering performance. To avoid these issues, it is crucial to dynamically orchestrate the distribution of the SoC power budget across the three domains based on their actual performance demands.

We propose SysScale, a new multi-domain power management technique to improve the energy efficiency of mobile SoCs. SysScale is based on three key ideas. First, SysScale introduces an accurate algorithm to predict the performance (e.g., bandwidth and latency) demands of the three SoC domains. Second, SysScale uses a new DVFS (dynamic voltage and frequency scaling) mechanism to distribute the SoC power to each domain according to the predicted performance demands. This mechanism is designed to minimize the significant latency overheads associated with applying DVFS across multiple domains. Third, in addition to using a global DVFS mechanism, SysScale uses domain-specialized techniques to optimize the energy efficiency of each domain at different operating points.

We implement SysScale on an Intel Skylake microprocessor for mobile devices and evaluate it using a wide variety of SPEC CPU2006, graphics (3DMark), and battery life workloads (e.g., video playback). On a 2-core Skylake, SysScale improves the performance of SPEC CPU2006 and 3DMark workloads by up to 16% and 8.9% (9.2% and 7.9% on average), respectively. For battery life workloads, which typically have fixed performance demands, SysScale reduces the average power consumption by up to 10.7% (8.5% on average), while meeting performance demands.

1. Introduction

A high-end mobile microprocessor is built as a system-on-chip (SoC) that integrates multiple components into a single chip. It typically has *three* main domains: *compute* (e.g., CPU cores, graphics engines), *IO* (e.g., display controller, image signal processing (ISP) engine), and *memory* (i.e., memory controller, memory interface, and DRAM) as illustrated in Fig. 1. A mobile SoC operates in a thermally-constrained environment, limited by what is known as thermal design power (TDP) [17, 19, 25, 34, 59, 66, 73, 83, 84, 86, 87, 93, 107]. To keep the

system running below a TDP, the SoC power-management-unit (PMU) employs a *power budget management* algorithm (PBM) to dynamically distribute the total power budget to each SoC domain [6, 18, 42, 57, 71, 80, 82, 84, 87, 111]. This allows each domain to operate within its allocated power budget. For instance, CPU cores and graphics engines in the compute domain share the same power budget. When a graphics-intensive workload is executed, the graphics engines consume most of the compute domain’s power budget. To keep the power consumption of the compute domain within its allocated power budget, PMU applies dynamic voltage and frequency scaling (DVFS) to 1) reduce the CPU cores’ power consumption and 2) increase the graphics engines’ performance [48, 51, 62, 77, 83, 84, 85, 87, 102, 109].

In this work, we demonstrate that the power budget the PBM allocates to the IO and memory domains is *inefficiently* managed, making the energy and performance of a high-end mobile SoC suboptimal. We make four key observations.

Observation 1. In a typical high-end mobile SoC, the power budget management algorithm assigns a *fixed* power budget to the IO and memory domains corresponding to the worst-case performance demands (bandwidth/latency) from the *IO interconnect* and the memory subsystem resources. However, we observe that common use cases of mobile systems have only modest demands relative to the worst-case. Unfortunately, these systems do not apply DVFS to the IO and memory domains based on the *actual* demands of the three domains, making these SoCs energy inefficient.

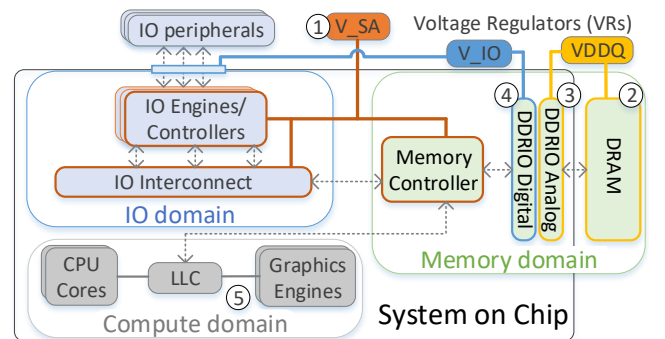


Figure 1: A modern mobile SoC (Intel Skylake [18]) with three domains (compute, IO, memory). Voltage regulators (VRs) are highlighted, e.g., IO engines/controllers, IO interconnect, and memory controller share the same VR, V_SA.

Observation 2. While mobile SoCs employ a power budget redistribution mechanism between components within a domain, such as between cores and graphics engines in the com-

pute domain [84, 85, 87], we observe that current power budget management algorithms *do not* support dynamic power redistribution *across* different domains. Therefore, when a domain’s power budget is underutilized, the remaining budget is wasted, making system performance suboptimal. This unused power budget could have been allocated to another domain (e.g., the compute domain) to increase performance. **Observation 3.** In modern mobile SoCs, we observe that multiple components in the IO and compute domains have widely-varying main memory bandwidth demands across different workloads. However, due to over-provisioning of IO and memory demands, SoC energy efficiency remains low while running many workloads, as we demonstrate in Sec. 3. **Observation 4.** Unoptimized DRAM configuration register values can significantly reduce the energy efficiency benefits of multi-domain DVFS (e.g., they provide 22% less power reduction than optimized values).

Unfortunately, there are *three* main challenges that make it difficult for existing high-end mobile systems to apply DVFS to the IO and memory domains based on demands of *multiple* domains. First, accurate prediction of 1) the actual bandwidth/latency demands of the multiple domains, and 2) the potential effect of DVFS on power/performance of the SoC, in the presence of multiple domains, is challenging. A modern high-end SoC integrates several components that share the IO interconnect and memory subsystem. Some of these components have strict quality of service (QoS) requirements [99] with respect to latency (e.g., isochronous traffic [5, 12, 103]) and bandwidth (e.g., display [9, 38, 94, 99]). Mispredicting a component’s actual demand can violate the QoS requirements and/or significantly degrade system performance. Second, the DVFS process of the IO and memory domains is a global system optimization. It requires monitoring the demands of the three SoC domains and subsequently configuring multiple components in the SoC to carry out the actual DVFS. Therefore, a power management transition flow for applying this global optimization can be computationally expensive. If it is not done correctly, the transition from one voltage/frequency operating point to another can degrade SoC performance by stalling the SoC domains. Third, the DVFS process should be holistic, efficient, and optimized to maximize power savings. For instance, previous works on memory subsystem DVFS [10, 11, 13, 14, 16, 20, 35, 58, 111] do not dynamically optimize the DRAM interface (i.e., DDRIO) *configuration registers* [52, 75, 79] and voltage during the DVFS process. Unoptimized DRAM configuration registers and voltage can significantly reduce, or even negate, the potential power/performance benefits of memory subsystem DVFS, as we show in this paper (Sec. 3).

Recent works in memory DVFS [11, 13, 14, 15, 16, 20, 35, 58, 111] for modern SoCs focus only on improving energy efficiency of a *single* domain (or limited components of two domains) and do not address all three challenges mentioned above. For example, MemDVFS [13] and MemScale [16] focus only on improving the energy efficiency of the main memory subsystem. CoScale [14] and other works [11, 20, 58] consider coordinating the DVFS of only CPU cores and the

main memory subsystem. To our knowledge, no previous work on SoC DVFS 1) coordinates and combines DVFS across three domains, or 2) optimizes the DRAM configuration registers [52, 75, 79] and voltage.

To enable more holistic power management in a mobile SoC and thereby to improve overall SoC power efficiency and performance, we propose *SysScale*, a new power management technique. SysScale is based on three **key ideas**. First, SysScale can accurately and dynamically predict the bandwidth/latency demand of multiple SoC domains by implementing new performance counters and utilizing existing system configuration registers. Second, SysScale uses a highly-efficient global DVFS mechanism to dynamically distribute the SoC power budget across all three domains, according to the predicted performance requirements. SysScale’s DVFS mechanism minimizes latency overheads by 1) performing DVFS simultaneously in all domains to overlap the DVFS latencies and 2) storing the required configuration registers in on-chip SRAM near each domain. Third, to maximize power savings, SysScale optimizes the energy efficiency of *each* domain at different DVFS operating points with domain-specific mechanisms. For instance, we optimize the energy efficiency of the DRAM interface by adding a dedicated scalable voltage supply and optimizing the configuration registers for each DVFS operating point.

This work makes the following **contributions**:

- To our knowledge, SysScale is the first work to enable coordinated and highly-efficient DVFS across all SoC domains to increase the energy efficiency of mobile SoCs. SysScale introduces the ability to redistribute the total power budget across all SoC domains according to the performance demands of each domain.
- We propose an effective algorithm to accurately predict the performance (e.g., bandwidth and latency) demands of the three SoC domains, utilizing newly-implemented performance counters and existing system configuration registers.
- We introduce a new global DVFS mechanism that minimizes the performance overhead of applying DVFS across multiple domains.
- We implement SysScale on the Intel Skylake SoC for mobile devices [2, 18, 97] and evaluate SysScale using a wide variety of workloads: SPEC CPU2006 [95], graphics (3DMark [100]), and battery life workloads for mobile devices [1] (e.g., web browsing, light gaming, video conferencing, and video playback). On a 2-core Skylake with a 4.5 W TDP, SysScale improves the performance of SPEC CPU2006 and 3DMark workloads by up to 16% and 8.9% (9.2% and 7.9% on average), respectively. For battery life workloads, which typically have fixed performance demands, SysScale reduces the average power consumption by up to 10.7% (8.5% on average), while meeting performance demands. As the TDP reduces, SysScale’s relative benefits significantly increase. For example, for a 3.5 W TDP system, SysScale enables up to 33% (19% on average) performance improvement on SPEC CPU2006 workloads.

2. Background

We provide a brief overview of a modern mobile SoC architecture and memory subsystem, with a focus on power consumption and DVFS.

2.1. Mobile SoC Architecture

Main Domains. A high-end mobile processor is a system-on-chip (SoC) that typically integrates three main domains into a single chip, as Fig. 1 shows: 1) compute (e.g., CPU cores and graphics engines), 2) IO (e.g., display controller, ISP engine, IO interconnect), and 3) memory (i.e., the memory controller, DRAM interface (DDRIO), and DRAM). The CPU cores, graphics engines, and IO controllers share the memory subsystem. Similarly, the IO controllers share the IO interconnect.

Clocks and Voltages. In a modern high-end mobile SoC, each one of the IO controllers/engines, IO interconnect, memory controller, and DDRIO typically have an independent clock. However, in current systems, there are three main voltage sources for the IO and memory domains. First, the IO controller, IO interconnect, and memory controller *share* the same voltage regulator, denoted as V_{SA} ¹ ① in Fig. 1. Second, DRAM ② and the analog part of the DRAM interface (DDRIO-analog) ③ share the same voltage regulator, known as VDDQ. Third, the digital part of the DRAM interface (DDRIO-digital) ④ typically shares the same voltage as the IO interfaces (e.g., display IO, ISP IO), denoted as V_{IO} in Fig. 1. The compute domain ⑤ typically has two voltage sources (not depicted in Fig. 1): 1) a voltage regulator that is shared between CPU cores and LLC and 2) a voltage regulator for the graphics engines [7, 8, 70, 74, 78, 91, 92, 97, 98].

2.2. Memory System Organization

We provide a general overview of the structure of a DRAM device. A set of DRAM chips placed together on a DRAM module comprises a rank. A DRAM chip is typically divided into multiple banks along with IO hardware (drivers and receivers) that enables access to the contents of the storage cells from outside the chip via the memory interface (DDRIO). Each bank includes peripheral logic to process commands and a grid of rows (wordlines) and columns (bitlines) of DRAM cells. A DRAM cell consists of a capacitor and an access transistor. A single bit of data is encoded by the charge level of the capacitor. For more detail, we refer the reader to prior works on DRAM organization and operation (e.g., [21, 22, 30, 31, 32, 49, 53, 54, 55, 56, 61, 68, 69, 88, 89, 90]).

2.3. Memory Power Consumption

The power consumption of a DRAM system comprises background power, operation power, and the memory controller power, which we briefly describe. We refer the reader to [10, 13, 21, 22] for a more detailed overview.

Background Power. A DRAM chip continuously consumes a certain amount of power in the background with or without memory accesses. The background power consumption has two main sources. First, the maintenance task that the pe-

ripheral circuitry performs to ensure signal integrity between the processor and DRAM chip. Second, periodic refresh operations that restore a leaky DRAM cell's charge level to ensure that a cell does not leak enough charge to cause a bit flip [22, 60, 61, 76].

Operation Power. Operation power is a DRAM chip's active power consumption when it executes memory commands for a memory access. It includes the DRAM array power, IO power, register power, and termination power.

DRAM array power is consumed by the core of the memory (e.g., banks, row/column decoders, and sense amplifiers). Thus, DRAM array power consumption correlates with DRAM access count (i.e., memory bandwidth consumption). The array draws a constant active-state power when a read, write, or precharge command is active. The IO power is consumed by input buffers, read/write latches, delay-locked loop (DLL), data interface drivers, and control logic that are used to transfer data from/to a DRAM chip. The register power is consumed by the input/output registers placed on clock and command/address DRAM interface lines, and their clock circuits (e.g., phase-locked loop (PLL)). The termination power includes the power used in terminating the DDRIO during active operation. Termination power depends on interface utilization and it is not directly frequency-dependent [13].

Memory Controller Power. Memory controller power is the combination of 1) the static power consumption, which is proportional to operational *voltage* and *temperature*, and 2) the dynamic power consumption, which is proportional to $voltage^2 \times frequency$ of the memory controller.

2.4. Memory DVFS

Dynamic voltage and frequency scaling (DVFS) is a technique that reduces the power consumption of an SoC component (e.g., a CPU core, a graphics engine) by reducing its voltage and frequency, potentially at the expense of performance [11, 13, 14, 15, 16, 20, 24, 26, 35, 58, 111]. While DVFS can allow a quadratic reduction in energy consumption and a cubic reduction in average power dissipation at the expense of a linear reduction in performance of the SoC component, the overall system energy consumption may *increase* due to longer execution time and higher utilization (i.e., less time spent in the idle power state) of other system components [24, 26]. Therefore, techniques such as race-to-sleep [85, 110] increase the frequency (and voltage) of an SoC component to reduce system-level energy. System architects use various metrics, such as the energy delay product (EDP)² to measure energy efficiency [23, 26], in a way that combines both energy and performance (delay).

The main objective of using DVFS in the memory subsystem of a mobile SoC is to improve energy efficiency (i.e., reduce EDP). To do so, memory DVFS reduces the background, operation, and memory controller power consumption by reducing the frequency and voltage of the memory subsystem. Current approaches scale the frequencies of the memory con-

¹SA stands for System Agent which houses the traditional Northbridge chip. SA contains several functionalities, such as the memory controller and the IO controllers/engines [106].

²Energy-delay product (EDP) [23] is a commonly used metric for quantifying a computing system's energy efficiency. The lower the EDP the better the energy efficiency.

troller, DDRIO (both analog and digital), and DRAM device, while the voltage is reduced only for the memory controller. Modifying the operating voltage of the DDRIO-analog and the DRAM device is *not* yet supported by commercially-available³ DRAM devices [10, 16, 46]. To maximize the energy savings while applying memory DVFS, we also concurrently apply DVFS to DDRIO-digital (④ in Fig. 1) and the IO interconnect. **Impact of Memory DVFS on the SoC.** Reducing the frequency of the memory subsystem affects the power consumption and performance of the SoC. *Performance* is affected because reducing the frequency 1) makes data bursts longer, 2) increases memory access time as it slows down both the memory controller and the DRAM interface, and 3) increases the queuing delays at the memory controller.

Power and energy consumption are also affected by reducing the memory subsystem frequency in four ways. First, background power reduces linearly. Second, memory controller power reduces approximately by a cubic factor due to the reduction of memory controller voltage (V_SA ① in Fig. 1) [11, 13, 14, 15, 16, 20, 24, 26, 35, 58, 111]. Third, lowering the DRAM operating frequency increases read, write, and termination energy linearly, because each access takes longer time. Fourth, due to the degradation in performance, the utilization of other components in the system can increase, which increases the overall SoC energy.

2.5. Memory Reference Code

Memory reference code (MRC [52, 75, 79]) is part of the BIOS code that manages system memory initialization. The purpose of MRC training is to 1) detect the DIMMs and their capabilities, 2) configure the configuration registers (CRs) of the memory controller (MC), DDRIO, and DIMMs, and 3) train the data and command interfaces (DDRIO) for correct operation and optimized performance as defined by the JEDEC standard [47].

MRC training is typically carried out with a *single* DRAM frequency. Therefore, the configuration register values of MC, DDRIO, and DIMMs are optimized for this particular DRAM frequency. When dynamically switching the memory subsystem between multiple frequencies (e.g., during DVFS), the configuration register values should be updated to the optimized values corresponding to the new DRAM frequency. Otherwise, the unoptimized configuration registers can degrade performance and negate potential benefits of DVFS, as we show in Sec. 3.

3. Motivation

To experimentally motivate building SysScale, we carry out an experiment on the Intel Broadwell processor [70], the previous generation of our target Skylake processor [18]. The goal of this experiment is to evaluate the potential benefits of employing DVFS across three SoC domains. We use multiple workloads from SPEC CPU2006 [95] and 3DMark [100], and a workload that exercises the peak memory bandwidth of DRAM [63]. We use two setups for our experiment, as we

show in Table 1: 1) To examine the performance of the processor without applying DVFS across multiple domains, we use a *baseline setup* in which we set the CPU core frequency to 1.2GHz and maintain the default voltage and frequency values of the other SoC components (e.g., DRAM, IO interconnect, and DDRIO digital). 2) To evaluate the benefits of applying DVFS across multiple domains, we use a *multi-domain DVFS setup* (MD-DVFS) in which the CPU cores have the same frequency as the first setup, but we reduce the frequency and voltage values of other SoC components in the IO and memory domains.

Table 1: Our two real experimental setups

Component	Baseline	MD-DVFS
DRAM frequency	1.6GHz	1.06GHz
IO Interconnect	0.8GHz	0.4GHz
Shared Voltage	V_SA	0.8·V_SA
DDRIO Digital	V_IO	0.85·V_IO
2 Cores (4 threads)	1.2GHz	1.2GHz

We attain the reduced-performance setup by changing *four* parameters. First, we reduce the DDR frequency by one bin⁴ (i.e., 1.06GHz). Doing so proportionally reduces the memory controller (MC) frequency that normally operates at half the DDR frequency. Since MC and IO interconnect share the same voltage (V_SA ① in Fig. 1), we also proportionally reduce the IO interconnect clock frequency to align it with the voltage levels of the IO and memory domains based on the voltage/frequency curves of both domains. Second, we reduce the shared voltage (V_SA ①) and DDRIO-digital voltage (V_IO ④) proportionally to the minimum functional voltage corresponding to the new frequencies in the IO and memory domains. Third, we maintain the CPU core (⑤) frequency and VDDQ (③) voltage unchanged across the two experimental setups. Fourth, to optimize the power/performance of the SoC in each setup, we configure the DRAM device, MC, and DDRIO with optimized MRC values for the selected DRAM frequency, as we explain in Sec. 2.5. Based on our evaluation of the SoC using the two setups, we make the following four key observations.

Observation 1. Current mobile systems have substantial energy inefficiency with respect to managing the voltage and frequency of the IO and memory domains.

We show in Fig. 2(a) the impact of the MD-DVFS setup on the average power consumption, energy consumption, performance, and energy-delay-product (EDP) compared to the baseline (without applying DVFS) using the three SPEC CPU2006 benchmarks. Fig. 2(a) also shows the performance of the MD-DVFS setup, compared to the baseline setup, when we increase the CPU core frequency from 1.2GHz to 1.3GHz.

We make three key observations from Fig. 2(a): First, the average power consumption of *all* three benchmarks reduces (by 10%–11%) with MD-DVFS. Second, while power consumption reduces in all evaluated workloads, in several workloads

³There are substantial challenges in operating a DRAM array at multiple voltages, which requires precisely tuned timing, transistor sizing, and voltage partitioning [46].

⁴DRAM devices support a few discrete frequency bins (normally only three). For example, LPDDR3 [45] supports only 1.6GHz, 1.06GHz, and 0.8GHz. The default bin for most systems is the highest frequency [18, 70].

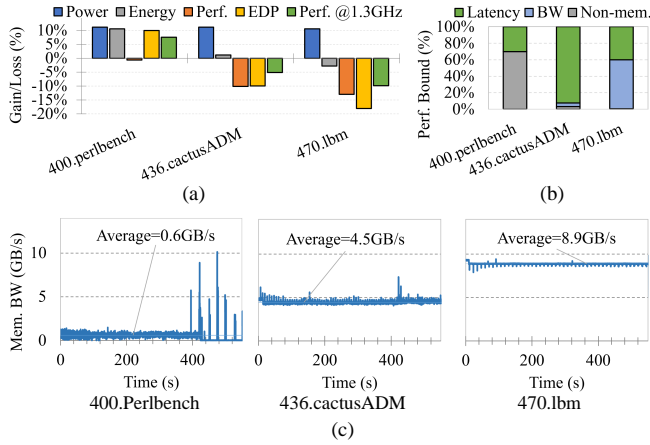


Figure 2: (a) Summary of the impact of MD-DVFS over the baseline on different metrics. (b) Bottleneck analysis of the three workloads, showing what fraction of the performance is bound by main memory latency, main memory bandwidth or non-main memory related events. (c) Memory bandwidth (BW) demand of the three benchmarks.

(e.g., cactusADM and lbm), there is a significant loss in performance ($>10\%$) as result of the reduction in frequency of the memory domain. Third, the effect of energy consumption varies widely across workloads. Workloads such as perlbench has reduced energy consumption by 11%. For workloads such as cactusADM and lbm, which experience a performance degradation with MD-DVFS, energy consumption only improves slightly or increases.

Fig. 2(b) shows a bottleneck analysis of the same three workloads evaluated in Fig. 2(a). We plot what fraction of the performance is bound by main memory latency, main memory bandwidth or non-main memory related events. We observe from Fig. 2(b) that cactusADM and lbm are heavily main memory bottlenecked, and the main bottleneck in cactusADM is main memory latency. Fig. 2(c) plots the memory bandwidth demand of the three benchmarks, showing that the memory bandwidth demand varies both over time and across workloads. Overall, the core-bound perlbench has low demand, but demand spikes at times, lbm has constant high demand, and cactusADM has moderate demand (but even then MD-DVFS hurts its performance by more than 10%).

We conclude that for workloads that are not memory latency or bandwidth bound, scaling down the voltage and frequency of the memory and IO domains can significantly reduce power and energy consumption with minimal effect on performance.

Observation 2. While mobile SoCs employ a power budget redistribution mechanism between SoC components within a domain, such as between CPU cores and the graphics engines in the compute domain [84, 85, 87], we observe that current mobile SoCs *do not* support dynamic power redistribution across different domains. Fig. 2(a) shows the performance impact of increasing the CPU cores’ frequency of the MD-DVFS setup from 1.2GHz to 1.3GHz when reassigning the saved average power budget from the IO and memory domains to the compute domain. Performance of perlbench improves sig-

nificantly by 8% over the baseline with MD-DVFS. Workloads that are not limited by compute bandwidth (e.g., cactusADM and lbm) do not experience performance improvement over the baseline with higher core frequency.

We conclude that 1) scaling down the voltage and frequency of the IO and memory domains when demands from these domains is low and 2) redistributing the saved power budget between domains can improve performance in workloads that are compute bound.

Observation 3. In modern mobile SoCs, we observe that multiple components in the IO (e.g., display controller, ISP engine) and compute (e.g., CPU cores, graphics engines) domains have widely-varying main memory bandwidth demands. Fig. 3(a) illustrates this by showing how the main memory bandwidth demand of three SPEC CPU2006 workloads and a 3DMark graphics workload varies over time. These workloads typically require different memory bandwidth over time.⁵ Fig. 3(b) shows the memory bandwidth demand of the display engine, ISP engine, and graphics engines (GFX) using different configurations and workloads. We observe that the display engine has widely-varying memory bandwidth demands depending on the display quality: HD display consumes approximately 17% of the peak memory bandwidth of a dual-channel LPDDR3 (25.6GB/s at 1.6GHz DRAM frequency), while a single 4K display (the highest supported display quality in our system) consumes 70% of the peak memory bandwidth.

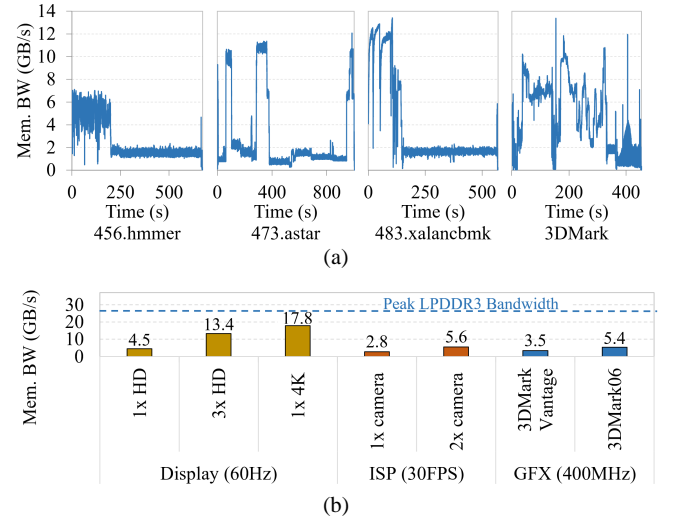


Figure 3: (a) Memory bandwidth (BW) demand over time for three SPEC CPU2006 benchmarks and the 3DMARK [100] graphics benchmark. (b) Average memory bandwidth demand of the display engine, ISP engine, and graphics engines (GFX) using different configurations and workloads.

We conclude that typical workloads have modest demands yet the SoC IO and memory demands are provisioned high, making existing mobile SoCs energy inefficient for typical workloads.

⁵The memory bandwidth demand was measured on the Intel Broadwell system [70] *without* scaling the voltage or frequency of any domain or component.

Observation 4. We observe that choosing optimized MRC values for the DRAM configuration registers is important for improving multi-domain DVFS energy efficiency. In Fig. 4, we show the impact of using *unoptimized* MRC values on the overall performance and power consumption of the MD-DVFS setup. We use a microbenchmark that was designed to exercise the peak memory bandwidth of DRAM (similar to STREAM [63, 67]). This helps us to isolate the impact of unoptimized MRC values on power/performance of the memory subsystem. We observe that unoptimized MRC values can greatly degrade both average power (by 22%) and performance (by 10%) compared to using optimized MRC values.

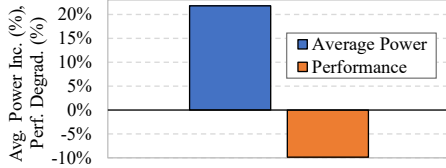


Figure 4: Impact of using unoptimized MRC values on power consumption and performance for a memory-bandwidth-intensive microbenchmark.

Based on our four key observations, we conclude that a *holistic power management approach* is needed to mitigate the power management inefficiencies in current mobile SoCs. This holistic approach should 1) redistribute the power budget of SoC domains based on the actual demands of a workload from each domain, 2) simultaneously scale the frequency/voltage of all domains, and 3) be optimized separately for each domain to minimize inefficiencies in the DVFS mechanism (e.g., by using optimized MRC values in the DRAM interface).

4. SysScale Architecture

We design SysScale with two design goals in mind: 1) reduce power consumption by dynamically scaling the voltage/frequency of all SoC domains based on performance demands from different domains, and 2) improve system throughput by redistributing the power budget across SoC domains based on performance demands.

SysScale achieves these two goals with *three* key components. The *first component* is a power management flow that is responsible for 1) scaling the multiple voltages and frequencies in the IO and memory domains and 2) reconfiguring the DRAM interface with optimized MRC values for the selected DRAM frequency.

The *second component* of SysScale is a *demand prediction mechanism* that uses both the system configuration and dedicated performance counters to predict the static and dynamic performance demands from the SoC domains. This component is important, as mispredicting the actual demand of the workload or an IO device places the system in an improper DVFS operating point, which can significantly degrade workload performance or violate quality of service requirements that the IO device might have.

The *third component* of SysScale is a *holistic power management algorithm* that is responsible for scaling the voltage and frequency of the IO interconnect and memory subsystem to meet the system’s dynamic performance demand. The power management algorithm uses the additional power budget saved as a result of DVFS in the IO and memory domains, to increase the performance of the compute domain (i.e., CPU cores and graphics engines). This leads to improved overall system throughput while maintaining the average system power within the thermal design power limit.

The three components of SysScale work together to orchestrate the voltage and frequency of each of the SoC domains and significantly reduce the overall energy consumption of a mobile SoC. We describe them in detail in the next three subsections.

4.1. Power Management Flow

SysScale power management flow, implemented by the PMU, is responsible for adjusting the frequencies and voltages of the IO interconnect and memory subsystem, as depicted in Fig. 5. First, SysScale’s *demand prediction mechanism* ❶ (described in Sec. 4.2) initiates a frequency change by determining new target frequencies/voltages for the SoC domains. If the demand prediction mechanism decides to increase (decrease) the IO and memory domain frequencies, then the flow increases (decreases) the voltages before ❷ (after ❷) the actual increase (decrease) of the PLL/DLL clock frequencies ❸. Subsequently, the flow blocks and drains the IO interconnect and the traffic from cores/graphics into the memory domain (i.e., LLC traffic to memory controller) ❹. To safely block the interconnect, all the outstanding requests are completed and new requests are not allowed to use the interconnect during this period. After all requests are completed, DRAM enters self-refresh mode ❺, and the flow loads the new optimized MRC values for the new DRAM frequency from on-chip SRAM into the memory-controller, DDRI/O and DRAM configuration registers ❻. Next, the SysScale flow sets the clocks to the new frequencies by re-locking both the phase-locked loops (PLLs) and delay-locked loops (DLLs) to the new IO interconnect and memory subsystem frequencies ❼. Finally, DRAM exits self-refresh mode ❽ and both the IO interconnect and the LLC traffic to the memory controller ❾ are released. This resumes SoC execution with the new frequencies/voltages for SoC domains and an optimized DRAM interface for the new DRAM frequency.

4.2. Demand Prediction Mechanism

The SysScale demand prediction mechanism uses peripheral configuration registers (e.g., the number of active displays or cameras) and new performance counters that we propose, to predict the performance demands of the SoC domains. We divide the performance demands into two categories: *static* and *dynamic* performance demands. We consider a performance demand as static if it is only related to system configuration (e.g., number of active displays or cameras). System configuration typically changes at the time-scale of tens of milliseconds as it is normally controlled by software (e.g., OS, drivers). Therefore, the PMU has enough time to re-

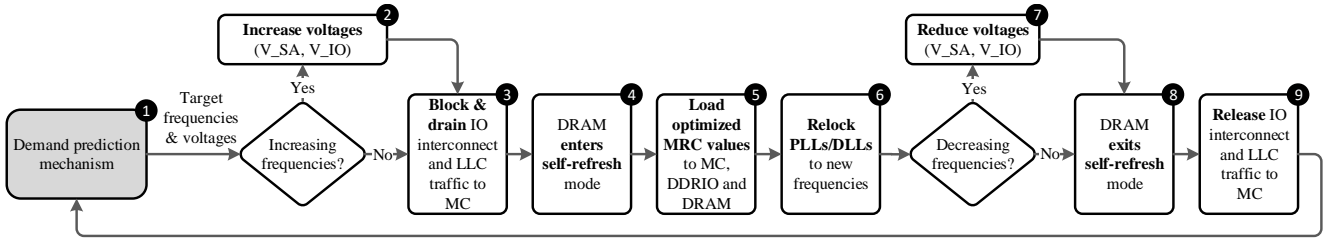


Figure 5: SysScale power management flow which carry out the DVFS of SoC domains.

spond to any configuration change that requires, for example, transition from a low- to high-bandwidth operating point of SysScale, without affecting system performance or quality of service. Such a transition can be completed within several *microseconds*. Dynamic performance demands are related to workload phase changes, which could happen much more frequently (e.g., hundreds of cycles). Next, we explain how the performance demand prediction is performed based on each demand category.

Static Performance Demand Estimation. To estimate static performance demand, SysScale maintains a *table* inside the firmware of the power-management unit (PMU) that maps every possible configuration of peripherals connected to the processor to IO and memory bandwidth/latency demand values. The firmware obtains the current configuration from control and status registers (CSRs) of these peripherals. For example, modern laptops support up to three display panels [104, 105]. When only one display panel is active (connected), it requires a certain bandwidth, but when three of the same display panel are connected, the bandwidth demand is nearly three times higher, as depicted in Fig. 3(b). In this case, the number of active displays and the resolution and refresh rate for each display are available in the CSRs of the display engine.

Our method enables accurate estimation of the static bandwidth demand based solely on the peripheral configuration as the bandwidth demand of a given peripheral configuration is known and is deterministic.

Dynamic Performance Demand Estimation. We categorize the bandwidth/latency demand of a workload as dynamic as it changes over time. For instance, SPEC CPU2006 workloads [95] demand different amounts of main memory bandwidth over time, as illustrated in Fig. 3. Similarly, mobile workloads have dynamically varying memory latency demands, as shown in [28].

We find that we are able to use existing performance counters to predict the dynamic bandwidth and latency demand when running a workload at a reduced frequency. We *select performance counters* for predicting dynamic IO and main memory bandwidth and latency demands using two steps: 1) among tens of internal processor performance counters, we build an initial selection based on intuition, and 2) we empirically prune our selection using an iterative process until the correlation between the performance counters and the workload’s performance degradation is closer to our target (e.g., >90% of our target). In the first stage, we choose an initial set of (15) internal performance counters related to memory

requests from the three domains of the SoC (i.e., compute, IO, and memory). Subsequently, we run a large number of representative mobile *workloads*⁶ with the two setups, baseline and multi-domain DVFS (MD-DVFS). We run the workloads on a real system with the same setup shown in Table 1, while selecting multiple DRAM frequencies for both the baseline and the MD-DVFS setup. We examine the performance of each run in addition to the values of the performance counters. Fig. 6 shows the actual versus the predicted performance impact of reducing DRAM frequency from a baseline frequency to a lower frequency used in MD-DVFS setup when running more than 1600 workloads. The figure also shows the correlation coefficient of the prediction when using the final list of performance counters. The final list of performance counters used in our prediction algorithm is as follows:

- **GFX_LLC_MISSES** counts the number of Last-Level-Cache (LLC) misses due to graphics engines (GFX). This performance counter is indicative of the bandwidth requirements of the graphics engines.
- **LLC_Occupancy_Tracer** provides the number of CPU requests that are waiting for data to return from the memory controller. This counter indicates whether the CPU cores are bandwidth limited.
- **LLC_STALLS** counts the number of stalls due to a busy LLC. This indicates that the workload is main memory latency limited.
- **IO_RPQ** (IO Read Pending Queue occupancy) counts the stalls due to busy IOs. This counter indicates that the workload is IO limited.

To *determine the threshold value* corresponding to each performance counter, we perform an offline phase that uses the results of the representative workloads (from the performance counter selection phase). For these runs, we set a bound on the performance degradation (e.g., 1%) when operating in MD-DVFS. We mark all the runs that have a performance degradation below this bound, and for the corresponding performance counter values, we calculate the *mean* (μ) and the *standard deviation* (σ). We set the threshold for each performance counter as $Threshold = \mu + \sigma$ [81].

Our *prediction algorithm* uses the performance counters to predict if the performance degradation of the running workload is less than a bound (e.g., 1%) when operating in MD-DVFS. To do so, the algorithm compares the value of each of the selected performance counters to its threshold.

⁶The workloads comprise of representative performance and office productivity workloads including SPEC06 [95], SYSmark [4], MobileMark [3], and computer graphics intensive workloads (e.g., 3DMARK [100]).

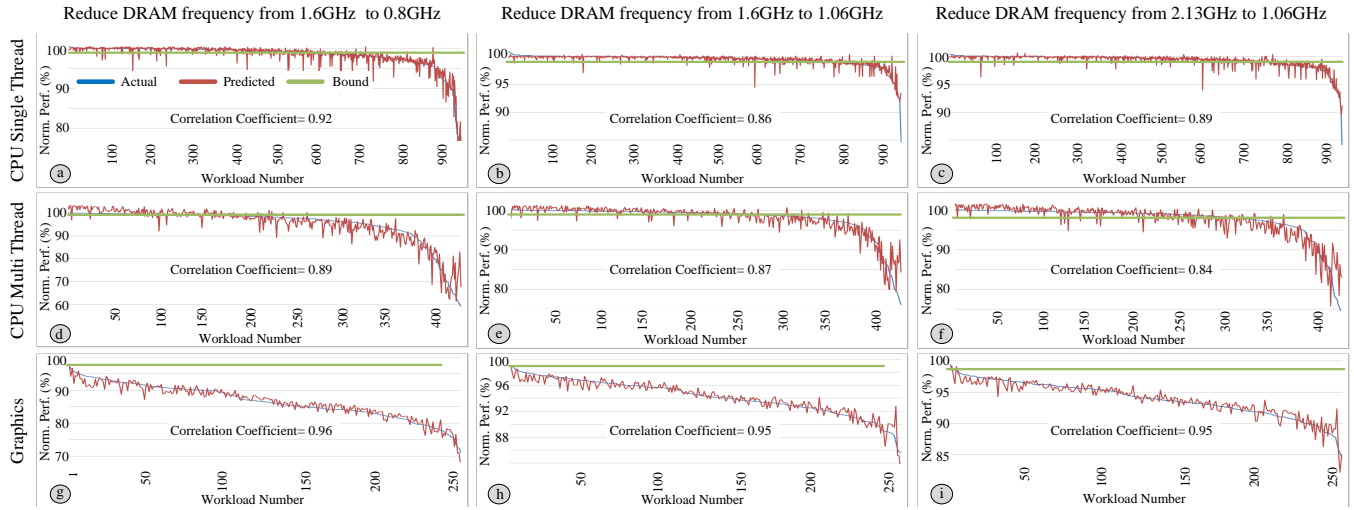


Figure 6: Actual vs. Predicted performance impact of reducing DRAM frequency in more than 1600 workloads. We examine three pairs (high/low) of DRAM frequencies: 1.6GHz/0.8GHz, 1.6GHz/1.06GHz and 2.13GHz/1.06GHz. Performance is normalized to that with the higher DRAM frequency. We also show the correlation coefficient between the Actual vs. Predicted performance impact values.

If the counter value is greater than its threshold, then the algorithm keeps the SoC at the current DVFS operating point; otherwise, the algorithm decides to move the IO and memory domains to the lower DVFS operating point.

Our prediction algorithm has an *accuracy* of 97.7%, 94.2%, and 98.8% for single threaded CPU, multi-threaded CPU, and graphics workloads, respectively (see Fig. 6). The prediction algorithm has no *false positive* predictions. This means that, there are no predictions where the algorithm decides to move the SoC to a lower DVFS operating point while the actual performance degradation is more than the bound.

4.3. Holistic Power Management Algorithm

SysScale implements a power distribution algorithm in the PMU firmware to manage multi-domain DVFS and redistribute power among the domains. PMU executes this algorithm periodically at a configurable time interval called *evaluation interval* (30ms by default). PMU samples the performance counters and CSRs (the configuration and status registers) multiple times in an evaluation interval (e.g, every 1ms) and uses the average value of each counter in the power distribution algorithm.

PMU switches the system between different performance levels (i.e., DVFS operating points) based on the predicted performance demand. Here we show the power distribution algorithm that switches the SoC between two operating points: high- and low-performance operating points. The system moves to the high-performance operating point if *any* of the following five conditions is satisfied.

1. The aggregated *static demand* requires higher memory bandwidth than a predefined threshold (STATIC_BW_THR).
2. The graphics engines are bandwidth limited (GFX_LLC_Misses > GFX_THR).
3. The CPU cores are bandwidth limited (LLC_Occupancy_Tracer > Core_THR).

4. Memory latency is a bottleneck (LLC_STALLS > LAT_THR)
5. IO latency is a bottleneck (IO_RPQ > IO_THR).

If *none* of these five conditions are true, PMU moves the system to the low-performance operating point. In the general case, where we have more than two SysScale operating points, the above algorithm decides between two adjacent operating points with dedicated thresholds.

When the SoC moves to the low-performance operating point, the PMU **reduces the power budgets of the IO and memory domains and increases the power budget of the compute domain** [18, 57, 84, 87].

The compute domain power budget management algorithm (PBM) distributes the received power budget between CPU cores and the graphics engines according to the characteristics of the workloads that run on these units. The additional power budget can increase the frequency of a thermally limited compute domain. PBM is designed to keep the average power consumption of the compute domain within the allocated power budget.

4.4. Interaction with CPU and Graphics DVFS

CPU cores and the graphics engines support independent DVFS mechanisms. DVFS states are known as P-states [24, 26, 37, 86]). The OS and the graphics driver normally initiate DVFS requests for the CPU cores and graphics engines, respectively. Each request is handled by the PBM algorithm inside the PMU firmware. PBM attempts to obey the DVFS request within the compute domain power budget constraint. If the request violates the power budget, then PBM demotes the request and places the requestor in a safe lower frequency that keeps the domain power within budget.

Compute domain DVFS mechanisms operate independently and are *not* directly tied to SysScale. SysScale interacts with compute domain DVFS mechanisms only via power budget management. When SysScale redistributes power budget from IO and memory domains to the compute domain, PBM

initiates DVFS in the compute domain to utilize the new budget and increase the domain’s frequencies. Similarly, when SysScale reduces the power budget of the compute domain (e.g., due to high IO or memory demand), PBM adjusts (if needed) the compute domain’s DVFS to keep the domain within the new power budget.

5. Implementation and Hardware Cost

SysScale requires the implementation of the three components inside the SoC.

First, SysScale requires the implementation of the four performance counters (as we describe in Sec. 4.2). The implementation of these counters in any SoC is straightforward.

Second, SysScale requires the implementation of hardware and firmware power management flow that enables the transition from one frequency/voltage operating point to another (as we summarize in Fig. 5). To implement this flow, two capabilities are needed: (1) the interconnect should support block and drain, which enables a seamless transition between frequencies; and (2) the flow should be able to configure the memory controller and DDRIO with the relevant MRC values for each frequency to support multiple optimized DRAM frequencies. The MRC values can be determined at reset time by performing MRC calculations for every supported frequency and saving the values inside an SRAM. The values are retrieved from the SRAM during the flow transition. To support MRC updates, we need to dedicate approximately 0.5KB of SRAM, which corresponds to less than 0.006% of Intel Skylake’s die area [18]. Additional capabilities, such as PLL/DLL re-locking and voltage adjustment, are supported in many SoCs. All of these capabilities exist inside the Skylake processor.

Finally, the actual demand prediction and power management algorithms are implemented inside the PMU firmware.⁷ The additional firmware code to support this flow is approximately 0.6KB —corresponding to less than 0.008% of Intel Skylake’s die area [18].

SysScale Transition Time Overhead. The actual latency of SysScale flow is less than 10 μ s. The latency has the following components: 1) voltage transitions (approximately $\pm 100mV$) of SA_V and IO_V voltage regulators (approximately 2 μ s with a voltage regulator slew rate of 50mV/ μ s), 2) draining IO interconnect request buffers (less than 1 μ s), 3) exiting DRAM self-refresh (less than 5 μ s with a fast training process), 4) loading the optimized DRAM configuration register from SRAM into configuration registers (less than 1 μ s), and 5) firmware latency and other flow overheads (less than 1 μ s).

6. Evaluation Methodology

We evaluate SysScale using real systems that employ Intel Broadwell [70] and Intel Skylake [18] SoCs. We use two distinct methodologies for 1) collecting motivational data and 2) evaluating SysScale. The reason is that we would like to

demonstrate the potential benefits of SysScale on the previous generation SoC (i.e., Broadwell) of our target Skylake SoC, before we implement it in the Skylake SoC.

Methodology for Collecting Motivational Data. To collect motivational data, we use a Broadwell-based system on which we emulate a crude version of SysScale’s static behavior, i.e., the multi-domain DVFS setup (MD-DVFS) that we use in Sec. 3. We use three steps to attain MD-DVFS setup. First, we reduce the DRAM frequency (to 1.06GHz) using the BIOS settings. Doing so boots the system with optimized MRC values for the DRAM interface (as we explain in Sec. 2.5), which exist in the BIOS. Second, since the memory controller and IO interconnect share the same voltage (V_SA ① in Fig. 1), we also proportionally reduce the IO interconnect clock frequency (to 0.4GHz) to align the voltage levels of the IO interconnect and memory controller. Third, we reduce the shared voltage and DDRIO voltage (i.e., V_SA and V_IO) by approximately 20% and 15%, respectively, which we determine based on the voltage/frequency curves of all components that share these voltages. To configure the voltage and frequency to the new values, we use the In-Target Probe (ITP) hardware debugger unit, which we explain below.

Methodology for Evaluating SysScale. We implement SysScale on the Intel Skylake SoC [18]. Table 2 shows the major system parameters. For our baseline measurements we disable SysScale on the same SoC.

Table 2: SoC and memory parameters

SoCs	M-5Y71 [40] Broadwell microarchitecture.
	M-6Y75 [41] Skylake microarchitecture.
Memory	CPU Core Base Frequency: 1.2GHz
	Graphics Engine Base Frequency: 300MHz
	L3 cache (LLC): 4MB.
	Thermal Design Point (TDP): 4.5W
	Process technology node: 14nm
	LPDDR3-1600MHz [45], non-ECC, dual-channel, 8GB capacity

In-Target Probe (ITP). ITP is a silicon debugger tool that connects to an Intel SoC through the JTAG port [108]. We use ITP 1) to set breakpoints at which the SoC halts when a specified event occurs and 2) configure the SoC’s control and status registers (CSRs) and model specific registers (MSRs) [36].

Power Measurements. We measure power consumption by using a National Instruments Data Acquisition (NI-DAQ) card (NI-PCIE-6376 [72]), whose sampling rate is up to 3.5 Mega-samples-per-second (MS/s). Differential cables transfer multiple signals from the power supply lines on the motherboard to the NI-DAQ card in the host computer that collects the power measurements. By using NI-DAQ, we measure power on up to 8 channels simultaneously. We connect each measurement channel to one voltage regulator of the SoC. The power measurement accuracy of the NI-PCIE-6376 is 99.94%.

Workloads. We evaluate SysScale with three classes of workloads that are widely used for evaluating mobile SoCs: 1) To evaluate CPU core performance, we use the SPEC

⁷A fully-hardware implementation is also possible. However, such power management flows are normally error-prone and require post-silicon tuning. As such, we choose to implement most of the flow within the power-management firmware (e.g., Pcode [24]).

CPU2006 benchmarks [95]. We use the SPEC CPU2006 benchmark score as the performance metric. 2) To evaluate computer graphics performance, we use the 3DMARK benchmarks [100]. We use frames per second (FPS) as the performance metric. 3) To evaluate the effect on battery life, we use a set of workloads that are typically used to evaluate the battery life of mobile devices such as, web browsing, light gaming, video conferencing, and video playback [1]. We use average power consumption as our battery life evaluation metric.

Methodology for Comparison to Prior Works. We compare SysScale to the two most relevant prior works, MemScale [16] and CoScale [14]. Unfortunately, 1) there is no real system available that implements these techniques, and 2) MemScale and CoScale techniques save *only* average power consumption without having the ability to redistribute excess power to other domains. In other words, they cannot boost the performance of the compute domain when applying DVFS to the memory domain. Therefore, to make our comparison useful, we *assume* that MemScale and CoScale can redistribute their saved power in the memory domain to increase the compute domain power budget (similar to SysScale). We call these improved versions of MemScale and CoScale, MemScale-Redist and CoScale-Redist, respectively (sometimes abbreviated as MemScale-R and CoScale-R).

We project the performance improvements of MemScale and CoScale compared to our baseline using three steps. First, we estimate the average power savings of MemScale and CoScale by using per-component measurements that we carry out on our Skylake system using power measurement tools that we describe in this section. Second, we build a performance/power model that maps an increase in the power budget of the compute domain to an increase in the CPU core or the graphics engine frequencies. For example, the model can show that a 100mW increase in compute power budget can lead to an increase in the core frequency by 100MHz. To build this model, we use performance and power measurements that we carry out on our Skylake system using multiple CPU core and graphics engine frequencies. Third, we use the *performance scalability*⁸ of the running workload with CPU frequency to project the actual performance improvement of the workloads with MemScale-Redist and CoScale-Redist.

7. Results

We present performance and average power benefits⁹ obtained with SysScale when it is implemented in an Intel Skylake SoC compared to the same SoC with SysScale disabled. We also compare SysScale to two prior works MemScale [16] and CoScale [14]. SysScale’s performance and average power benefits are *measured* results, while the results for MemScale-Redist and CoScale-Redist are *projected*, as we explain in

⁸We define performance scalability of a workload with respect to CPU frequency as the performance improvement the workload experiences with unit increase in frequency.

⁹Energy efficiency, i.e., energy-delay-product (EDP) [23], improves proportionally to performance or average power consumption, since SysScale improves either performance within a fixed power budget or average power consumption within a fixed performance requirement.

Sec. 6. We evaluate three workload categories: CPU (Sec. 7.1), graphics (Sec. 7.2), and battery life workloads (Sec. 7.3). We also analyze sensitivity to different SoC thermal-design-power (TDP) levels and DRAM frequencies (Sec. 7.4).

7.1. Evaluation of CPU Workloads

Fig. 7 reports the performance improvements of MemScale-R, CoScale-R, and SysScale over our baseline system. We make four key observations.

First, SysScale improves real system performance by 9.2% on average. This result is significant as it is obtained on a real system.

Second, SysScale provides $5.4\times$ and $2.4\times$ the performance improvement of MemScale-R and CoScale-R, respectively. This is because 1) SysScale is holistic, taking into account all SoC domains (and components in each domain), whereas MemScale and CoScale consider only memory DVFS and coordinated (CPU and memory subsystem) DVFS, respectively, and 2) MemScale and CoScale do not dynamically optimize the DRAM interface configuration registers and voltage after each DVFS change. Also recall that the performance improvement of MemScale-R and CoScale-R is projected based on the estimated average power reduction of each technique (Sec. 6), whereas SysScale improvements are real measured improvements.

Third, the performance benefit of SysScale correlates with the *performance scalability* of the running workload with CPU frequency. Highly-scalable workloads (i.e., those bottlenecked by CPU cores, such as *416.gamess* and *444.namd*) have the highest performance gains. In contrast, workloads that are heavily bottlenecked by main memory, such as *410.bwaves* and *433.milc*, have almost no performance gain. We note that a highly scalable workload (e.g., *416.gamess*) benefits in two ways from SysScale: 1) because the workload is not bottlenecked by main memory, reducing memory frequency reduces power consumption without affecting performance, 2) because the workload is bottlenecked by CPU performance, redistributing the power budget from memory and IO to the CPU increases the workload’s performance.

Fourth, if a workload has different execution phases bottlenecked by different SoC domains (e.g., compute versus memory), then SysScale dynamically adjusts the frequencies of the multiple domains to improve system performance. For example, *473.astar* has execution phases of up to several seconds of low memory bandwidth demand (e.g., $\sim 1\text{GB/s}$) and high memory bandwidth demand (about $\sim 10\text{GB/s}$), as illustrated in Fig. 3(a), and SysScale significantly improves its performance by 13%.

We conclude that SysScale significantly improves CPU core performance by holistically applying DVFS to SoC domains based on dynamic performance demands and dynamically redistributing the excess power budget between domains, while carefully optimizing the memory interface during DVFS transitions.

7.2. Evaluation of Graphics Workloads

Typically, the performance of a graphics workload is highly scalable with the graphics engine frequency. When running

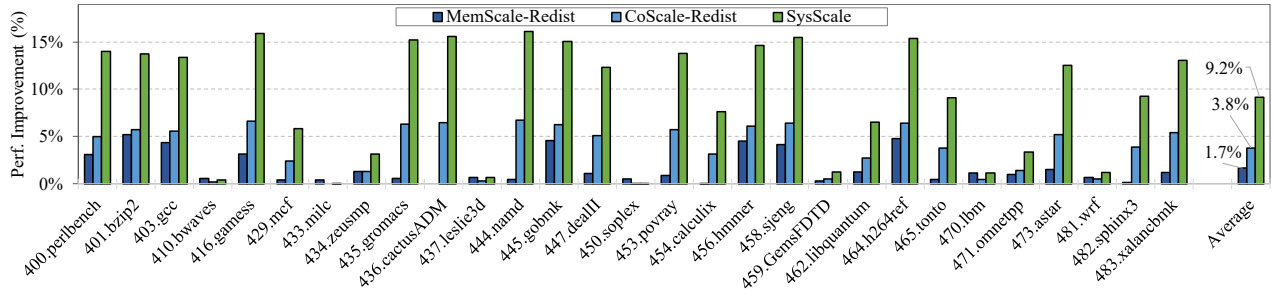


Figure 7: Performance improvement of MemScale-Redist, CoScale-Redist, and SysScale on SPEC CPU2006 workloads.

graphics workloads, the power budget management algorithm (PBM [57, 83]) of the PMU normally allocates only 10% to 20% of the compute domain power budget to the CPU cores, while the graphics engines consumes the rest of the power budget [82, 84, 87]. For a mobile system, while running a graphics workload, the CPU cores normally run at the most energy efficient frequency P_n [26] (i.e., the maximum possible frequency at the minimum functional voltage (V_{min})). Moreover, at a very low TDP (e.g., 3.5 W and 4.5 W), the effective CPU frequency is reduced below P_n by using hardware duty cycling (HDC¹⁰) [37].

Fig. 8 shows the performance improvement of MemScale-R, CoScale-R, and SysScale, compared to our baseline when running three different 3DMark [100] graphics workloads. We make two key observations. First, SysScale improves the system performance of 3DMark06, 3DMark11, and 3DMark Vantage by 8.9%, 6.7% and 8.1%, respectively. SysScale improves performance because it boosts the graphics engines, frequency by redistributing the power budget across three SoC domains.

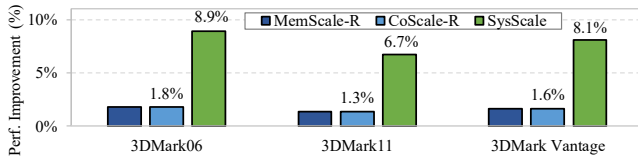


Figure 8: Performance improvement of MemScale-R, CoScale-R, and SysScale on computer graphics workloads.

Second, SysScale provides approximately $5\times$ the performance improvement of MemScale-R and CoScale-R. MemScale-R and CoScale-R have similar performance improvements because the average power savings of these two techniques is identical. The reason is that, in these workloads, the CPU cores run at the lowest possible frequency, and therefore CoScale (which applies DVFS to CPU cores in addition to the memory subsystem) cannot further scale down the CPU frequency.

We conclude that the saved power budget from IO and memory domains can be used to raise the frequency of the graphics engines and improve graphics performance.

¹⁰Hardware Duty Cycling, also known as SoC Duty Cycling, implements coarse grained duty cycling by using C-states with power gating in contrast to T-states that use clock gating [37].

7.3. Evaluation of Battery Life Workloads

Unlike CPU and graphics workloads that always benefit from higher performance, the battery life workloads have two characteristics: 1) their performance demands are fixed; for example, in video playback running at 60 frames per second, each frame needs to be processed and displayed on the display panel within 16.67 milliseconds, and 2) they have long idle phases where the system enters into idle power states (C-states [24, 26, 27, 101]). We note that according to our measurements, the active state (i.e., C0 power state) residency of these workloads is between 10%-40%, while the SoC is in idle state (e.g., C2, C6, C7, or C8 [24, 26, 27, 101]) during the remaining execution time of the workload. In the C0 power state, typically the compute domain (i.e., CPU cores and graphics engines) operates in the lowest possible frequencies, while in all other power states, the compute domain is idle (i.e., clock- or power-gated).

Fig. 9 shows the SoC average power reduction when running four representative battery life workloads [1], web-browsing, light-gaming, video conferencing, and video-playback, with a single HD display panel (e.g., the laptop display) is active. We make three key observations.

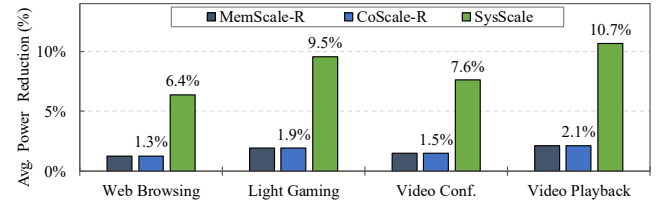


Figure 9: Average power reduction of MemScale-R, CoScale-R, and SysScale on representative battery life workloads.

First, SysScale reduces average power consumption of web browsing, light gaming, video conferencing, and video playback workloads by 6.4%, 9.5%, 7.6%, and 10.7%, respectively, on our real system.

Second, SysScale provides approximately $5\times$ the power reduction of MemScale-R and CoScale-R. MemScale-R and CoScale-R provide similar average power reduction benefits to each other, since in battery life workloads, the CPU cores run at the lowest possible frequency, and therefore CoScale cannot further reduce the CPU cores, frequency.

Third, SysScale provides significant power savings for battery life workloads for all SoC power states in which DRAM

is active (i.e., not in self-refresh). For example, in the video playback workload, the SoC transitions between three power states during the processing of each video frame: C0, C2, and C8 with residencies of 10%, 5%, and 85%, respectively (not graphed). DRAM is active only in C0 and C2 power states, while in C8, DRAM is in self-refresh. Therefore, SysScale applies DVFS to the IO and memory domains only during C0 and C2 states in the video playback workload.

We conclude that for battery life workloads, which have fixed performance requirements, applying SysScale as a holistic DVFS mechanism for all three SoC domains significantly reduces the average SoC power consumption.

7.4. Sensitivity Analysis

In this section, we investigate the effect of system parameters on SysScale performance and power gains.

SoC Thermal Design Point (TDP). The evaluation results presented in Sections 7.1, 7.2, and 7.3 are for an SoC (i.e., Skylake M-6Y75 [41]) with a TDP of 4.5W, as we show in Table 2. This SoC has a configurable TDP ranging from 3.5W up to 7W. The Skylake architecture itself can scale up to a TDP of 91W [39] for a high-end desktop.

Fig. 10 shows the average performance improvement on SPEC CPU2006 workloads when running SysScale on systems with different TDPs. Violin plots show the distribution and density of performance improvement points across different workloads compared to each TDP’s baseline. We make two key observations. First, at a constrained TDP of 3.5W, SysScale improves performance by up to 33% (19.1% on average). Second, as TDP increases, SysScale’s performance benefit reduces. This is because power becomes ample and there is less need to redistribute it across domains.

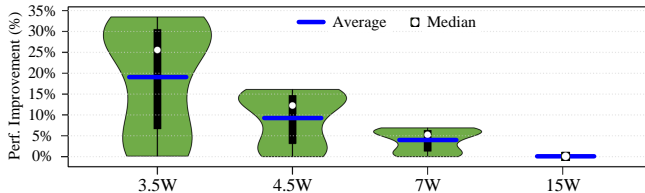


Figure 10: SysScale performance benefit vs. TDP on SPEC CPU2006 workloads.

We also evaluate battery life workloads and found that SysScale’s average power savings are not affected by the TDP. In these workloads, the compute domain typically operates at the lowest possible frequencies for the CPU cores and graphics engines (i.e., the P_n P-state [24, 26]) regardless of the TDP of the SoC.

We conclude that SysScale significantly improves the performance of especially TDP-constrained SoCs while it improves battery life (energy consumption) across the entire TDP range (e.g., 3.5W–91W [39, 41]).

More DRAM Frequencies. In our previous experimental results, we use a system with a base DRAM frequency of 1.6GHz. We also evaluate a system with different DRAM devices and frequencies. Based on our evaluation of SysScale

with different DRAM types and frequencies, we make three key observations.

First, SysScale supports other DRAM devices with different DRAM frequencies. The average power savings (freed-up power budget that is dynamically utilized to boost the compute domain) obtained by SysScale when scaling DDR4’s operating point from 1.86GHz to 1.33GHz is approximately 7% lower than that when scaling LPDDR3 operating point from 1.6GHz down to 1.06GHz.

Second, SysScale supports multiple operating points for the same DRAM device. For example, we could use 0.8GHz as the lowest DVFS operating point of the memory subsystem. However, we observe that the 0.8GHz operating point is not energy efficient when compared to 1.06GHz due to two main reasons. 1) The system agent voltage (V_{SA}) already reaches the minimum functional voltage (V_{min}) when scaling DDR frequency to 1.06GHz. Therefore, reducing the DRAM frequency beyond 1.06GHz does not provide significant benefits as the V_{SA} voltage remains the same. 2) When comparing Figures 6(a,d,g) and Figures 6(b,e,h), we observe that the average performance degradation of the traces when reducing the DRAM frequency from 1.6GHz down to 0.8GHz is $2 \times - 3 \times$ higher than that when reducing it from 1.6GHz down to 1.06GHz. As such, we implement only two operating points (i.e., 1.6GHz and 1.06GHz) in our real system.

Third, a **finer-grained** DVFS of the three SoC domains can increase the benefits of SysScale. For instance, supporting additional DDR frequencies between the 1.6GHz to 0.8GHz range can increase the performance and/or the average power of workloads with memory demand that fall in between existing frequencies. Unfortunately, current commercial DDR devices support only few frequencies (e.g., LPDDR3 supports only 1.6GHz, 1.06GHz, and 0.8GHz).

We conclude that the benefits of SysScale can be higher with more control over the DRAM frequencies. To achieve this, we recommend that the **DRAM vendors** enable finer-grained control on DDR frequency in future generations to help increase the energy efficiency of modern mobile SoCs.

8. Related Work

To our knowledge, this is the first work to 1) enable coordinated and highly-efficient DVFS across all SoC domains to increase the energy efficiency of mobile SoCs and 2) provide multi-domain DVFS results from real a modern mobile SoC (Intel Skylake) running real mobile system benchmarks.

Compute Domain DVFS. Many prior works propose DVFS for the compute domain (i.e., the CPU cores and graphics engines) [24, 26, 29, 33, 35, 37, 42, 43, 44, 48, 50, 64, 65, 86, 96, 111]. DVFS operating points for the compute domain, known as P-states [24, 26, 37, 86], are typically managed by the OS and the graphics drivers for the CPU cores and graphics engines, respectively. These mechanisms optimize voltage and frequency assuming a fixed power budget for the compute domain. Unlike SysScale, they do not enable optimization and distribution of power budget across different SoC domains and thus provide limited energy efficiency benefits in mobile SoCs.

Memory DVFS and Coordinated DVFS. Recent works in memory DVFS [10, 11, 13, 14, 15, 16, 20, 35, 58, 111] for modern SoCs focus only on improving energy efficiency of the memory domain (e.g., MemDVFS [13], Voltron [10], and MemScale [16]) or limited components of two domains (e.g., CoScale [14] and other works [11, 20, 58]).

In Sections 7.1, 7.2, and 7.3, we qualitatively and quantitatively compare SysScale to two closely related prior works, MemScale [16] and CoScale [14]. We show that SysScale significantly increases system performance and reduces energy consumption compared to the two mechanisms due to two reasons: 1) SysScale is holistic, taking into account all SoC domains (and components in each domain), whereas MemScale and CoScale consider only memory DVFS and/or compute domain DVFS and 2) MemScale and CoScale do not dynamically optimize DRAM interface configuration registers and voltage after each DVFS change.

Current memory DVFS approaches [10, 11, 13, 14, 16, 20, 35, 58, 111] all have one or more of three main drawbacks that make them inefficient for modern mobile SoCs. First, many of these prior works are not directly optimized for modern mobile SoC architectures that integrate *many* components into the three SoC domains in a very thermally constrained environment. Many works target server architectures, which have different structure, constraints, and workloads than in mobile systems. Second, past works scale *only* one or two system domains at a time. Such limited scaling provides small benefits for mobile SoCs (as we show in Sections 3 and 7) and, when accounting for the relatively high DVFS transition costs, the benefits diminish even further. Third, past works do not dynamically optimize the DRAM interface configuration registers [52, 75, 79], which degrades performance and may negate potential gains (as we show in Sections 3 and 7).

9. Conclusion

We propose SysScale, the first work to enable coordinated and highly-efficient DVFS across all SoC domains to increase the energy efficiency of mobile SoCs. SysScale introduces the ability to optimize and efficiently redistribute the total power budget across all SoC domains according to the performance demands of each domain. SysScale is implemented in the Intel Skylake SoC for mobile devices. We show that it significantly improves the performance of real CPU and graphics workloads (by up to 16% and 8.9%, respectively, for 4.5 W TDP) and reduces the average power consumption of battery life workloads (by up to 10.7%) on a real Intel Skylake system. We conclude that SysScale is an effective approach to balance power consumption and performance demands across all SoC domains in a sophisticated heterogeneous mobile SoC to improve energy efficiency and performance.

Acknowledgments

We thank the anonymous reviewers of ISCA 2020 for feedback and the SAFARI group members for feedback and the stimulating intellectual environment they provide.

References

- [1] Anandtech, "The Microsoft Surface Pro (2017) Review: Evaluation," online, accessed March 2018, Mar 2019, <https://bit.ly/2VEdhCg>.
- [2] I. Anati, D. Blythe, J. Doweck, H. Jiang *et al.*, "Inside 6th-Generation Intel® Core™: New microarchitecture code named Skylake," in *Hot Chips*, 2016.
- [3] BAPCo, "MobileMark 2014," online, accessed May 2018, Mar 2019, <https://bapco.com/products/mobilemark-2018>.
- [4] BAPCo, "SYSmark 2014," online, accessed May 2018, Mar 2019, <https://bapco.com/products/sysmark-2014-se/>.
- [5] L. Benini, G. De Micheli, and T. T. Ye, *Networks on Chips*. Morgan Kaufmann Burlington, 2006.
- [6] P. S. Bhojwani, J. D. Lee, and R. N. Mahapatra, "SAPP: Scalable and Adaptable Peak Power Management in NoCs," in *ISLPED*, 2007.
- [7] T. Burd, N. Beck, S. White, M. Paraschou, N. Kalyanasundharam, G. Donley, A. Smith, L. Hewitt, and S. Naffziger, "Zeppelin: An SoC for Multi-chip Architectures," *JSSC*, 2019.
- [8] E. A. Burton, G. Schrom, F. Paillet, J. Douglas, W. J. Lambert, K. Radhakrishnan, and M. J. Hill, "FIVR - Fully Integrated Voltage Regulators on 4th Generation Intel® Core™ SoCs," in *APEC*, 2014.
- [9] J. A. Castellano, *Handbook of Display Technology*. Elsevier, 2012.
- [10] K. K. Chang, A. G. Yaglikci, S. Ghose, A. Agrawal, N. Chatterjee, A. Kashyap, D. Lee, M. O'Connor, H. Hassan, and O. Mutlu, "Understanding Reduced-voltage Operation in Modern DRAM Devices: Experimental Characterization, Analysis, and Mechanisms," *SIGMETRICS*, 2017.
- [11] J. Chen and L. K. John, "Predictive Coordination of Multiple On-chip Resources for Chip Multiprocessors," in *ICS*, 2011.
- [12] W. J. Dally, "Interconnect-centric computing," in *HPCA*, 2007.
- [13] H. David, C. Fallin, E. Gorbato, U. R. Hanebutte, and O. Mutlu, "Memory Power Management via Dynamic Voltage/Frequency Scaling," in *ICAC*, 2011.
- [14] Q. Deng, D. Meisner, A. Bhattacharjee, T. F. Wenisch, and R. Bianchini, "CoScale: Coordinating CPU and Memory System DVFS in Server Systems," in *MICRO*, 2012.
- [15] Q. Deng, D. Meisner, A. Bhattacharjee, T. F. Wenisch, and R. Bianchini, "MultiScale: memory system DVFS with multiple memory controllers," in *ISLPED*, 2012.
- [16] Q. Deng, D. Meisner, L. Ramos, T. F. Wenisch, and R. Bianchini, "MemScale: Active Low-power Modes for Main Memory," in *ASPLOS*, 2011.
- [17] J. Donald and M. Martonosi, "Techniques for Multicore Thermal Management: Classification and New Exploration," *ISCA*, vol. 34, no. 2, pp. 78–88, 2006.
- [18] J. Doweck, W.-F. Kao, A. K.-y. Lu, J. Mandelblat *et al.*, "Inside 6th-Generation Intel Core: New Microarchitecture Code-Named Skylake," *IEEE Micro*, 2017.
- [19] H. Esmaeilzadeh, E. Blem, R. S. Amant, K. Sankaralingam, and D. Burger, "Dark Silicon and the End of Multicore Scaling," in *ISCA*, 2011.
- [20] W. Felter, K. Rajamani, T. Keller, and C. Rusu, "A Performance Conserving Approach for Reducing Peak Power Consumption in Server Systems," in *ICS*, 2005.
- [21] S. Ghose, T. Li, N. Hajinazar, D. S. Cali, and O. Mutlu, "Demystifying Complex Workload-DRAM Interactions: An Experimental Study," *SIGMETRICS*, 2019.
- [22] S. Ghose, A. G. Yaglikci, R. Gupta, D. Lee, K. Kudrolli, W. X. Liu, H. Hassan, K. K. Chang, N. Chatterjee, A. Agrawal *et al.*, "What your DRAM power models are not telling you: Lessons from a detailed experimental study," *SIGMETRICS*, 2018.
- [23] R. Gonzalez and M. Horowitz, "Energy Dissipation in General Purpose Microprocessors," *JSSC*, 1996.
- [24] C. Gough, I. Steiner, and W. Saunders, "CPU Power Management," in *Energy Efficient Servers: Blueprints for Data Center Optimization*. Springer, 2015.
- [25] P. E. Gronowski, W. J. Bowhill, R. P. Preston, M. K. Gowan, and R. L. Allmon, "High-Performance Microprocessor Design," *JSSC*, 1998.
- [26] J. Haj-Yahya, A. Mendelson, Y. B. Asher, and A. Chattopadhyay, "Power Management of Modern Processors," in *Energy Efficient High Performance Processors*. Springer, 2018.
- [27] J. Haj-Yahya, Y. Sazeides, M. Alser, E. Rotem, and O. Mutlu, "Techniques for Reducing the Connected-Standby Energy Consumption of Mobile Devices," in *HPCA*, 2020.
- [28] J. Haj-Yahya, A. Yasin, and Y. Ben-Asher, "DOEE: Dynamic Optimization Framework for Better Energy Efficiency," in *HiPC*, 2015.
- [29] V. Hanumaiah and S. Vrudhula, "Energy-Efficient Operation of Multicore Processors by DVFS, Task Migration, and Active Cooling," *IEEE TC*, 2012.
- [30] H. Hassan, M. Patel, J. S. Kim, A. G. Yaglikci, N. Vijaykumar, N. M. Ghiasi, S. Ghose, and O. Mutlu, "CROW: A Low-Cost Substrate for Improving DRAM Performance, Energy Efficiency, and Reliability," in *ISCA*, 2019.
- [31] H. Hassan, G. Pekhimenko, N. Vijaykumar, V. Seshadri, D. Lee, O. Ergin, and O. Mutlu, "ChargeCache: Reducing DRAM Latency by Exploiting Row Access Locality," in *HPCA*, 2016.
- [32] H. Hassan, N. Vijaykumar, S. Khan, S. Ghose, K. Chang, G. Pekhimenko, D. Lee, O. Ergin, and O. Mutlu, "SoftMC: A Flexible and Practical Open-source Infrastructure for Enabling Experimental DRAM Studies," in *HPCA*, 2017.
- [33] S. Herbert and D. Marculescu, "Analysis of Dynamic Voltage/Frequency Scaling in Chip-Multiprocessors," in *ISLPED*, 2007.
- [34] W. Huang, M. R. Stan, K. Sankaranarayanan, R. J. Ribando, and K. Skadron, "Many-core Design from a Thermal Perspective," in *DAC*, 2008.
- [35] C. Imes, H. Zhang, K. Zhao, and H. Hoffmann, "Handing DVFS to hardware: Using Power Capping to Control Software Performance," *Technical Report*, 2018.
- [36] Intel, "In Target Probe (ITP)," online, accessed July 2019, <https://bit.ly/2qHSMbm>.
- [37] Intel, "Intel 64 and IA-32 Architectures Software Developer's Manual Volume 3A, 3B, and 3C," online, accessed July 2019, <https://intel.ly/2S5i9OY>.

- [38] Intel, "Open Source Graphics," <https://bit.ly/3aJXy9q>, 2016.
- [39] Intel, "Intel Core i7-6700K Processor," online, accessed June 2019, <https://intel.ly/2KzUUly>.
- [40] Intel, "Intel Core M-5Y71 Processor," online, accessed June 2019, <https://intel.ly/3eSNhLi>.
- [41] Intel, "Intel Core M-6Y75 Processor," online, accessed June 2019, <https://intel.ly/2xdbyL4>.
- [42] C. Isci, A. Buyuktosunoglu, C.-Y. Cher, P. Bose, and M. Martonosi, "An Analysis of Efficient Multi-core Global Power Management Policies: Maximizing Performance for a Given Power Budget," in *MICRO*, 2006.
- [43] C. Isci, A. Buyuktosunoglu, and M. Martonosi, "Long-term Workload Phases: Duration Predictions and Applications to DVFS," *IEEE Micro*, 2005.
- [44] C. Isci, G. Contreras, and M. Martonosi, "Live, Runtime Phase Monitoring and Prediction on Real Systems with Application to Dynamic Power Management," in *MICRO*, 2006.
- [45] JEDEC, "JESD79-3-1.35 V DDR3L-800, DDR3L-1066, DDR3L-1333, DDR3L-1600, and DDR3L-1866," *JEDEC Std., JESD79-3-1A*, vol. 1, 2013.
- [46] JEDEC, "Low Power DRAM Evolution," online, 2018, <https://bit.ly/2rrqxq>.
- [47] JEDEC, "JEDEC Standard," online, Apr 2020, <https://www.jedec.org/sites/default/files/docs/JESD212.pdf>.
- [48] W. Kim, M. S. Gupta, G.-Y. Wei, and D. Brooks, "System Level Analysis of Fast, Per-core DVFS Using On-chip Switching Regulators," in *HPCA*, 2008.
- [49] Y. Kim, V. Seshadri, D. Lee, J. Liu, and O. Mutlu, "A Case for Exploiting Subarray-level Parallelism (SALP) in DRAM," in *ISCA*, 2012.
- [50] T. Kolpe, A. Zhai, and S. S. Sapatnekar, "Enabling Improved Power Management in Multicore Processors Through Clustered DVFS," in *DATE*, IEEE, 2011.
- [51] T. Komoda, S. Hayashi, T. Nakada, S. Miwa, and H. Nakamura, "Power capping of CPU-GPU heterogeneous systems through coordinating DVFS and task mapping," in *ICCD*, 2013.
- [52] D. Lai, A. Chen, and Y. Li, "MRC Training vs. Board Defect," in *IMPACT*, 2014.
- [53] D. Lee, S. Ghose, G. Pekhimenko, S. Khan, and O. Mutlu, "Simultaneous Multi-Layer Access: Improving 3D-Stacked Memory Bandwidth at Low Cost," in *TACO*, 2016.
- [54] D. Lee, Y. Kim, G. Pekhimenko, S. Khan, V. Seshadri, K. Chang, and O. Mutlu, "Adaptive-latency DRAM: Optimizing DRAM Timing for the Common-case," in *HPCA*, 2015.
- [55] D. Lee, Y. Kim, V. Seshadri, J. Liu, L. Subramanian, and O. Mutlu, "Tiered-Latency DRAM: A Low Latency and Low Cost DRAM Architecture," in *HPCA*, 2013.
- [56] D. Lee, L. Subramanian, R. Ausavarungnirun, J. Choi, and O. Mutlu, "Decoupled Direct Memory Access: Isolating CPU and IO Traffic by Leveraging a Dual-Data-Port DRAM," in *PACT*, 2015.
- [57] O. Lempel, "2nd Generation Intel® Core Processor Family: Intel® Core i7, i5 and i3," in *Hot Chips*, 2011.
- [58] X. Li, R. Gupta, S. V. Adv, and Y. Zhou, "Cross-Component Energy Management: Joint Adaptation of Processor and Memory," *TACO*, 2007.
- [59] Y. Li, B. Lee, D. Brooks, Z. Hu, and K. Skadron, "CMP Design Space Exploration Subject to Physical Constraints," in *HPCA*, 2006.
- [60] J. Liu, B. Jaiyen, Y. Kim, C. Wilkerson, and O. Mutlu, "An Experimental Study of Data Retention Behavior in Modern DRAM Devices: Implications for Retention Time Profiling mechanisms," *ISCA*, 2013.
- [61] J. Liu, B. Jaiyen, R. Veras, and O. Mutlu, "RAIDR: Retention-Aware Intelligent DRAM Refresh," in *ISCA*, 2012.
- [62] A. Malik, B. Lin, G. Memik, P. Dinda, and R. P. Dick, "User-driven Frequency Scaling," *CAL*, 2006.
- [63] J. D. McCalpin *et al.*, "Memory Bandwidth and Machine Balance in Current High Performance Computers," *TCCA*, 1995.
- [64] A. Miyoshi, C. Lefurgy, E. Van Hensbergen, R. Rajamony, and R. Rajkumar, "Critical Power Slope: Understanding the Runtime Effects of Frequency Scaling," in *ICS*, 2002.
- [65] C. Möbius, W. Dargie, and A. Schill, "Power Consumption Estimation Models for Processors, Virtual Machines, and Servers," *TPDS*, 2013.
- [66] M. Monchiero, R. Canal, and A. González, "Design Space Exploration for Multi-core Architectures: A Power/Performance/Thermal View," in *ICS*, 2006.
- [67] T. Moscibroda and O. Mutlu, "Memory Performance Attacks: Denial of Memory Service in Multi-core Systems," in *USENIX Security*, 2007.
- [68] O. Mutlu and T. Moscibroda, "Stall-Time Fair Memory Access Scheduling for Chip Multiprocessors," in *MICRO*, 2007.
- [69] O. Mutlu and T. Moscibroda, "Parallelism-Aware Batch Scheduling: Enabling High-performance And Fair Shared Memory Controllers," in *ISCA*, 2008.
- [70] A. Nalamalpu, N. Kurd, A. Deval, C. Mozak, J. Douglas, A. Khanna, F. Paillet, G. Schrom, and B. Phelps, "Broadwell: A Family of IA 14nm Processors," in *VLSI Circuits*, 2015.
- [71] R. Nathuji, K. Schwan, A. Somani, and Y. Joshi, "VPM Tokens: Virtual Machine-aware Power Budgeting in Datacenters," *Cluster computing*, 2009.
- [72] National Instruments, "NI-DAQ PCIe-6376," online, accessed 2019, <http://www.ni.com/pdf/manuals/377387c.pdf>.
- [73] A. Naveh, E. Rotem, A. Mendelson, S. Gochman, R. Chabukswar, K. Krishnan, and A. Kumar, "Power and Thermal Management in the Intel Core Duo Processor," *ITJ*, 2006.
- [74] V. P. Nikolskiy, V. V. Stegailov, and V. S. Vecher, "Efficiency of the Tegra K1 and X1 SoC for Classical Molecular Dynamics," in *HPCS*, 2016.
- [75] P. Nukala, S. More, C. Mozak, D. Morgan, and R. Loop, "Enhanced Data Transfer for Multi-Loaded Source Synchronous Signal Groups," in *EPEPS*, 2018.
- [76] M. Patel, J. S. Kim, and O. Mutlu, "The Reach Profiler (REAPER) Enabling the Mitigation of DRAM Retention Failures via Profiling at Aggressive Conditions," *ISCA*, 2017.
- [77] A. Pathania, Q. Jiao, A. Prakash, and T. Mitra, "Integrated CPU-GPU Power Management for 3D Mobile Games," in *DAC*, 2014.
- [78] Qualcomm Technologies, "Qualcomm Snapdragon 410E (APQ8016E) Processor Device Specification," online, 2018, <https://bit.ly/2VByEEo>.
- [79] B. Querbach, R. Khanna, D. Blankenbeckler, Y. Zhang *et al.*, "A Reusable BIST with Software Assisted Repair Technology for Improved Memory and IO Debug, Validation and Test Time," in *ITC*, 2014.
- [80] P. Ranganathan, P. Leech, D. Irwin, and J. Chase, "Ensemble-level Power Management for Dense Blade Servers," *ISCA*, 2006.
- [81] C. Reimann, P. Filzmoser, and R. G. Garrett, "Background and Threshold: Critical Comparison of Methods of Determination," *STOTEN*, 2005.
- [82] E. Rotem, A. Naveh, D. Rajwan, A. Ananthakrishnan, and E. Weissmann, "Power Management Architecture of the 2nd Generation Intel® Core Microarchitecture, Formerly Codenamed Sandy Bridge," in *Hot Chips*, 2011.
- [83] E. Rotem, "Intel Architecture, Code Name Skylake Deep Dive: A New Architecture to Manage Power Performance and Energy Efficiency," in *IDF*, 2015.
- [84] E. Rotem, R. Ginosar, A. Mendelson, and U. C. Weiser, "Power and Thermal Constraints of Modern System-on-a-Chip Computer," in *THERMINIC*, 2013.
- [85] E. Rotem, R. Ginosar, C. Weiser, and A. Mendelson, "Energy Aware Race to Halt: A Down to EARTH Approach for Platform Energy Management," *IEEE CAL*, 2014.
- [86] E. Rotem, A. Mendelson, R. Ginosar, and U. Weiser, "Multiple Clock and Voltage Domains for Chip Multi Processors," in *MICRO*, 2009.
- [87] E. Rotem, A. Naveh, A. Ananthakrishnan, E. Weissmann, and D. Rajwan, "Power-Management Architecture of the Intel Microarchitecture Code-Named Sandy Bridge," *IEEE Micro*, 2012.
- [88] V. Seshadri, Y. Kim, C. Fallin, D. Lee, R. Ausavarungnirun, G. Pekhimenko, Y. Luo, O. Mutlu, P. B. Gibbons, M. A. Kozuch, and T. Mowry, "RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization," in *MICRO*, 2013.
- [89] V. Seshadri, D. Lee, T. Mullins, H. Hassan, A. Boroumand, J. Kim, M. A. Kozuch, O. Mutlu, P. B. Gibbons, and T. C. Mowry, "Ambit: In-memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology," in *MICRO*, 2017.
- [90] V. Seshadri and O. Mutlu, "In-DRAM Bulk Bitwise Execution Engine," *arXiv preprint arXiv:1905.09822*, 2019.
- [91] T. Singh, S. Rangarajan, D. John, C. Henrion, S. Southard, H. McIntyre, A. Novak, S. Kosonocky, R. Jotwani, A. Schaefer *et al.*, "3.2 Zen: A Next-Generation High-Performance x86 Core," in *ISSCC*, 2017.
- [92] T. Singh, A. Schaefer, S. Rangarajan, D. John, C. Henrion, R. Schreiber, M. Rodriguez, S. Kosonocky, S. Naffziger, and A. Novak, "Zen: An Energy-Efficient High-Performance x86 Core," *JSSC*, 2018.
- [93] K. Skadron, M. R. Stan, K. Sankaranarayanan, W. Huang, S. Velusamy, and D. Tarjan, "Temperature-aware Microarchitecture: Modeling and Implementation," *TACO*, 2004.
- [94] Y. Song, O. Alavoine, and B. Lin, "A Self-aware Resource Management Framework for Heterogeneous Multicore SoCs with Diverse QoS Targets," *TACO*, 2019.
- [95] Standard Performance Evaluation Corporation, "SPEC," online, accessed March 2018, www.spec.org.
- [96] V. Sundriyal, M. Sosonkina, B. Westheimer, and M. Gordon, "Core and Uncore Joint Frequency Scaling Strategy," *JCC*, 2018.
- [97] S. M. Tam, H. Muljono, M. Huang, S. Iyer, K. Royneogi, N. Satti, R. Qureshi, W. Chen, T. Wang, H. Hsieh *et al.*, "Skylake-SP: A 14nm 28-Core Xeon® Processor," in *ISSCC*, 2018.
- [98] Z. Toprak-Deniz, M. Sperling, J. Bulzacchelli, G. Still, R. Kruse, S. Kim, D. Boerstler, T. Gloekler, R. Robertazzi, K. Stawiasz *et al.*, "5.2 Distributed System of Digitally Controlled Microregulators Enabling Per-Core DVFS for the POWER8 TM Microprocessor," in *ISSCC*, 2014.
- [99] H. Usui, L. Subramanian, K. K.-W. Chang, and O. Mutlu, "DASH: Deadline-aware High-performance Memory Scheduler for Heterogeneous Systems with Hardware Accelerators," *TACO*, 2016.
- [100] Vantage, "3DMARK," online, accessed March 2018, <http://www.futuremark.com/benchmarks/3dmarkvantage>.
- [101] Vantage, "Advanced Configuration and Power Interface (ACPI) specification," online, accessed January 2019, <http://www.acpi.info>.
- [102] P.-H. Wang, Y.-M. Chen, C.-L. Yang, and Y.-J. Cheng, "A Predictive Shutdown Technique for GPU Shader Processors," *CAL*, 2009.
- [103] W.-D. Weber, J. Chou, I. Swarbrick, and D. Wingard, "A Quality-of-Service Mechanism for Interconnection Networks in System-on-Chips," in *DATE*, 2005.
- [104] Wikipedia, "Apple MacBook Air," online, accessed June 2019, https://en.wikipedia.org/wiki/MacBook_Air.
- [105] Wikipedia, "Microsoft Surface," online, accessed June 2019, https://en.wikipedia.org/wiki/Microsoft_Surface.
- [106] Wikipedia, "Skylake Microarchitecture," online, accessed June 2019, <https://bit.ly/2NG8Hz2>.
- [107] Wikipedia, "Thermal Design Power," online, accessed September 2019, https://en.wikipedia.org/wiki/Thermal_design_power.
- [108] M. Williams, "Low Pin-Count Debug Interfaces for Multi-Device Systems," *ARM's Serial Wire Debug white paper*, 2009.
- [109] Q. Wu, P. Juang, M. Martonosi, and D. W. Clark, "Voltage and Frequency Control with Adaptive Reaction Time in Multiple-clock-domain Processors," in *HPCA*, 2005.
- [110] H. Zhang, P. V. Rengasamy, S. Zhao, N. C. Nachiappan, A. Sivasubramanian *et al.*, "Race-to-Sleep+ Content Caching+ Display Caching: A Recipe for Energy-Efficient Video Streaming on Handhelds," in *MICRO*, 2017.
- [111] H. Zhang and H. Hoffmann, "Maximizing Performance Under a Power Cap: A Comparison of Hardware, Software, and Hybrid Techniques," in *ASPLOS*, 2016.