

# Many-Core Token-Based Adaptive Power Gating

Andrew B. Kahng<sup>†‡</sup>, Seokhyeong Kang<sup>‡</sup>, Tajana Rosing<sup>†‡</sup> and Richard Strong<sup>†</sup>

<sup>†</sup>CSE and <sup>‡</sup>ECE Departments, University of California at San Diego

La Jolla, CA 92093-0404 USA

abk@ucsd.edu, shkang@vlsicad.ucsd.edu, tajana@ucsd.edu, rstrong@eng.ucsd.edu

**Abstract**—Among power dissipation components, leakage power has become more dominant with each successive technology node. Leakage energy waste can be reduced by power gating. In this work, we extend Token-Based Adaptive Power Gating (TAP), a technique to power gate an *actively* executing core during memory accesses, to many-core CMPs. TAP works by tracking every system memory request and its estimated time of arrival so that a core may power gate itself without performance or energy loss. Previous work on TAP [11] shows several benefits compared to earlier state-of-the-art techniques [10], including zero performance hit and  $2.58\times$  average energy savings for out-of-order cores. We show that TAP can adapt to increasing memory contention by increasing power-gated time by  $3.69\times$  compared to a low memory-pressure case. We also scale TAP to many-core architectures with a distributed wake-up controller that is capable of supporting staggered wake-ups and able to power gate each core for 99.07% of the time achieved by a non-scalable centralized scheme.

## I. INTRODUCTION

During every cycle that a core is on, even when stalled, leakage power is consumed via gate leakage, gate-induced drain leakage, junction leakage, and subthreshold leakage. A core may stall quite often if it is intensely accessing the memory subsystem, as every time a thread makes a memory request that misses in the L1 cache, the core is subjected to a variable access latency. This variable latency often translates into a core stall during which no forward thread progress occurs and energy is wasted. For a 32nm out-of-order EV6 core, stall energy can be up to 39.1% of total energy consumption for the Spec2006 benchmarks [17].

*Power gating* is a technique that drastically reduces leakage power by cutting off the current path from supply to ground through introduction of a transistor switch between them. At one end of the spectrum, functional unit power gating reduces power consumption of unused core functional units [25] with wake-up latencies of several nanoseconds. At the other end, entire cores may be power gated and woken up, with latencies of several tens of microseconds to account for saving and restoring all core state from memory [22]. An intermediate mechanism, Memory Access Power Gating [10], provides the ability to power gate an entire core, wake up a power-gated core in about 10ns, and maintain the core's architectural and cache state.

In this work, we extend *Token-Based Adaptive Power Gating* (TAP) [11] to many-core architectures. TAP deterministically applies power gating during core stalls which are caused

by the variable latency of requests to the memory subsystem. TAP achieves this by providing the capability to track every ongoing memory request and the expected response time for each memory access that misses in the L1 cache. An expected lower bound on latency is sent to each core's PPGS by modifying the cache controllers to send a token on any miss where the token includes an estimate of the access latency of a next-level memory hit. The result is that TAP can support power gating with no performance loss.

Previous work on TAP [11] assumes a centralized wake-up controller (WUC) which is a critical limiter to the scalability of the system. During each power gating action, the core must communicate with the WUC to ensure that it can safely use its wake-up mode without violating voltage noise constraints. We replace the WUC with a distributed wake-up scheme that assigns cores to predetermined, recurring wake-up slots that enforce safe wake-up modes across the CMP. Such a scheme is shown to operate nearly as efficiently as an ideal WUC. Our key contributions are the following:

- We introduce a distributed wake-up scheme to control core wake-up modes in many-core designs that sacrifices only 0.93% of power-gated time.
- We demonstrate that TAP can adapt to an increase in memory contention by increasing power-gated time by  $3.69\times$  as the number of threads increases from 1 to 32.
- We design and implement a staggered wake-up scheme capable of reducing wake-up latency by up to 58.2%.

## II. RELATED WORK

Power gating has been studied at both circuit and architectural levels. Circuit-level papers typically analyze different circuit techniques aimed at reducing wake-up latency, efficiently retaining logic states, minimizing ground bounce, and achieving resilience to process variation. Microarchitectural works typically examine questions related to use of different power-gating modes, what to power gate, predicting when to power gate, and control algorithms to avoid energy penalties from poor power-gating decisions. The following briefly reviews representative works in these two areas.

In the realm of circuit innovation, the pioneering work of Horiguchi et al. [8] has been followed by many works on fundamental circuit design issues related to power gating, including switch-cell sizing, data-retention methods, physical-implementation methodologies, and mode-transition noise analysis and reduction. The recent survey of Shin et al. [23] gives an excellent summary of the history and highlights of power-gating techniques. Agarwal et al. [4] and Singh et al. [24] examine multiple sleep modes that feature different wake-up overheads and leakage power savings. To minimize

Authors are listed alphabetically by last name. Principal contributors, to whom correspondence should be addressed: Seokhyeong Kang (shkang@vlsicad.ucsd.edu) and Richard Strong (rstrong@eng.ucsd.edu).

ground bounce during mode transition, Kim et al. [12] control turn-on voltage ( $V_{GS}$ ), which makes sleep transistors turn on in a non-uniform stepwise manner. Chowdhury et al. [6] propose a tri-mode power-gating technique using PMOS switches in parallel with NMOS footer switches. Finally, Zhang et al. [26] propose a multi-mode power-gating technique using three NMOS switches with different sizes and threshold voltages. Using various combinations of the three switches, they can provide multiple power-gating modes with different leakage savings.

In the microarchitectural arena, Hu et al. [9] propose use of power gating to reduce functional unit leakage power when applications underutilize their functional units. [9] also develops equations to estimate the break-even points for power gating an out-of-order superscalar processor. Lungu et al. [18] show that in many cases, the predictor of [9] can lead to increased energy consumption. Madan et al. [19] extend the idea of Lungu et al. to the core level and propose a “guard mechanism” that reduces harmful use of power gating.

Power-gating technology is also readily visible in leading commercial products. The recent Nehalem architecture employs power gating at the core level to reduce leakage power on idle cores, but  $100ms$  is required to wake up a core [13], [16]. AMD [22] has improved this power-gating technique by optimizing the wake-up sequence to skip built-in self tests (BIST) and restoration of cache state; this results in wake-up times as short as  $75\mu s$ . In today’s systems, the OS typically power gates the cores in the idle loop, missing out on power gating long memory accesses.

Li et al. [14] consider a scheme that directs core DVFS behavior based on signals from the L2-cache and estimates of instruction-level parallelism to reduce energy consumption on memory bound applications, but does not consider power gating. Our early work, Memory Access Power Gating [10], provides a mechanism to power gate an entire core, wake up a power-gated core in about  $10ns$ , and maintain the core’s architectural and cache state. It uses a combination of a programmable power-gating switch (PPGS), state-retention cells, and source biasing to enable the core to efficiently enter and exit a power-gated state. The mechanism is shown to be capable of reducing leakage power during memory accesses. A similar work, Memory Access Aware Power Gating for MPSoCs [15], examines the potential to power gate an in-order core while monitoring a single memory bus and estimating memory latencies. Token-Based Adaptive Power Gating (TAP) [11] extends memory access power gating to out-of-order execution and considers using staggered wake-up for cores, but avoids questions about the scalability of TAP to many-core architectures. By contrast, our present work addresses the scalability of TAP to many-core designs.

### III. SYSTEM DESIGN

A memory access power-gating controller must provide three functions. First, the controller should be able to predict the expected duration of core stalls. Second, the controller must assign a core’s programmable power-gating switch (PPGS) [10] a wake-up mode that does not violate supply voltage noise constraints of the system when waking up a core. Last, the controller must retain essential core architectural and performance-related state. Together, these functions allow for energy savings and minimal performance hit without violating

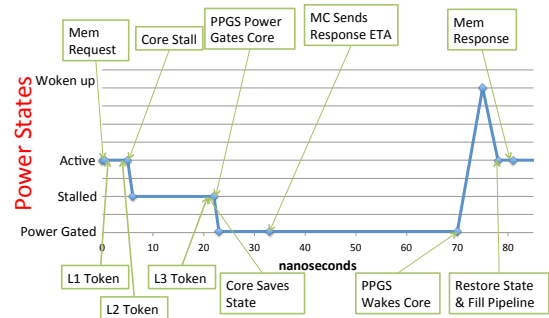


Fig. 1. Diagram depicting the core going through various power states (Power Gated, Stalled, Active, Woken up) as the PPGS power gates the core on a memory access.

voltage noise constraints. The rest of this section describes how our system provides these three functions.

**TAP: Token-Based Latency Monitoring.** TAP informs each PPGS about expected memory latency by modifying the cache controllers to send tokens on cache misses that include an estimate of the lower-bound access latency of a next-level memory hit derived from Table I and a time stamp of creation<sup>1</sup>. The controllers send the tokens to the PPGS of the core that requested the memory access. Once the PPGS receives the token, it looks at the lower-bound latency to satisfy the request and power gates the core if the core is both stalled and idle long enough to save energy. Should the core receive more than one token for simultaneous memory requests, it tracks each expected response separately and schedules the resumption of core execution to satisfy the earliest response. If a token is delayed in the memory subsystem by a controller or queue, the PPGS can compare its arrival time with its generation time stamp and previous tokens to determine whether the token should be ignored.

Whenever a memory request misses all the way to the memory controller, the response latency experiences a significant amount of variability. This variability is caused by the complexity of DRAM memory [27], which includes bank queues, availability of the data in the row-buffer, writing wrong address row-buffers, accessing the column in the row-buffer, and channel contention between banks. TAP adapts to memory variability by adding a special token. As soon as the last-level cache experiences a miss, a token is sent to the requesting core’s PPGS with an estimated completion time of *UNKNOWN*. This is a directive to the PPGS to start power gating its core immediately and to expect one additional token with the ETA of the memory response. Once the memory controller submits the memory access to one of the banks and determines whether the access is a row-buffer hit or miss, it sends the second *ETA* token to the core’s PPGS with the ETA of the response assuming that there is no memory channel contention. The PPGS receives the second token before the response and schedules core wake-up for the appropriate time.

Figure 1 shows a timing-accurate diagram of a PPGS power gating the core in response to messages from the TAP technique. At time  $0ns$ , a memory request occurs that will miss in the cache hierarchy and cause a memory access. The PPGS then receives tokens for the L1, L2 and L3 misses. Just after receiving the L2 token, the core stalls due to a dependency.

<sup>1</sup>A cache controller sitting on the core side of a shared NUCA cache would require a per-bank lower-bound access latency.

After the L3 token is received, the PPGS decides to power gate the core and saves all core state. The core is then power gated and the memory controller (MC) sends a updated ETA for the memory response. At  $70ns$ , the PPGS begins waking up the core. At  $78ns$ , the core state is restored and the pipeline is restarted. The memory response comes back at  $81ns$  and the core resumes execution as if nothing happened.

### Distributed Staggered Wake Up.

Previous work [10] designed its Wake-up Controller (WUC) for the worst case when all cores wake up simultaneously. However, wake-up latency is significantly reduced when we *stagger* the wake-up sequence so that two cores wake up at slightly different times (e.g., offset by  $1ns$ ). In such a case, stagger reduces the worst-case peak current and voltage noise that a core may experience. We now consider how to integrate stagger into a many-core design; analysis of the benefit of stagger is given in Section IV-B.

To use staggered wake up, a controller must give wake-up modes and times to each core. For a large multi-core system (e.g., 64 cores), a given core's PPGS may not tolerate the latency to communicate with a centralized WUC due to propagation and queuing delay across the chip. However, we observe that non-adjacent cores do not significantly affect core wake-up latency, and with proper guardband, only adjacent cores (eight cores at most) need be considered.

This observation motivates a distributed design which assigns each core to a recurring wake-up slot. In this scheme, each core is given a recurring slot at which it can start waking up. To avoid any performance hit, a core should select a slot before the deadline to start waking up. The average wake-up delay for a core is defined by two degrees of freedom: the number of unique wake-up slots,  $\eta$ , and the stagger between two adjacent wake-up slots,  $\psi$ . The worst-case reduction in power-gated time occurs when a core predicts that it would need to wake up  $\epsilon$  seconds before its assigned slot such that  $\epsilon < \psi$ , causing the core to wake up  $\eta * \psi - \epsilon$  seconds earlier. Given a core that wakes up at any time, uniformly at random, the average expected reduction in power-gated time is  $\frac{\eta * \psi}{2}$ .

Values of  $\eta$  and  $\psi$  should be chosen to maximize a core's power-gated time and minimize the safe wake-up latency of the core. Given  $\eta$ , wake-up slots should be assigned to cores to minimize the number of adjacent woken-up cores. Figure 2 shows three ways of assigning wake-up modes to cores such that the number of adjacent woken-up cores is minimized with a preference given to cores waking up simultaneously in the diagonal position. An increase in  $\eta$  and  $\psi$  reduces the maximum number of simultaneous core wake-ups and the minimum safe wake-up latency. At the same time, increasing these two parameters increases average expected reduction in power-gated time. According to our simulations, system energy savings are maximized when  $\eta$  and  $\psi$  equal 5 (no simultaneous wake-ups) and  $0.9ns$ , respectively. This setting results in a  $10.3ns$  minimal wake-up latency and  $2.25ns \pm 1.30ns$  average reduction in power-gated time per core power-gating opportunity, compared to  $9.6ns$  with  $0.9ns$  stagger for an ideal centralized WUC, when simulated on GEM5 [5] with the Spec2006 benchmarks.

**Core State Retention and Restoration.** To avoid losing core state that is required for correct and efficient execution, essential sequential and SRAM cells must be retained. We use the

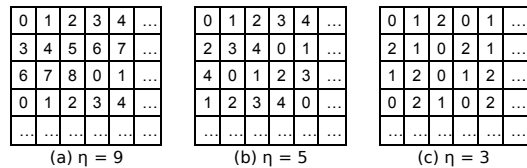


Fig. 2. Wake-up slot assignments with different number of slots ( $\eta$ ).

technique from [10] which replaces a subset of sequential cells with live-slave retention flip-flops [7] which can be triggered to retain their logical values before a power-gating action at a cost of 20% increase in area and power versus a normal flip-flop. Only those sequential cells comprising the architectural registers necessary to refill the pipeline are selected, which results in a 3.4% area overhead for the processor. SRAM cells' state is retained through source biasing [21] in which the supply voltage is reduced to 50% of nominal so that SRAM leakage is reduced, but logical state is maintained. This technique enables saving contents of L1 caches, TLBs, branch predictor state, physical registers, etc. To provide supply power during power gating, a separate non-collapsible voltage domain provides power to the retention flip-flops and SRAM cells. Thus, as the power is gated from combinational logic and non-essential sequential cells, the separate voltage rail provides power to maintain core state. The overhead from multi-power domains and separate voltage rails already exists for power gating cores today. Additional cycles are required for the power gating and wake-up sequence, and to disable/enable the clock, trigger data retention, refill the pipeline, and de-assert/assert the clamps. We model the power-down and wake-up sequence as in [10]. Figure 3 shows an in-order core implementation for the power gating and restoration with retention flip-flops.

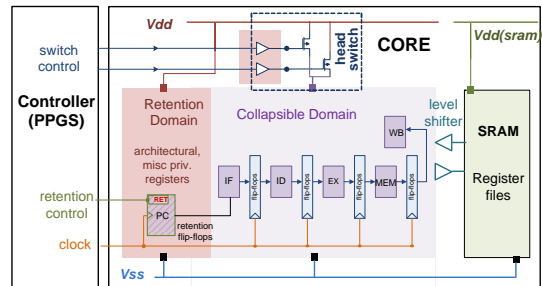


Fig. 3. Interface for power gating and data retention.

## IV. RESULTS

Table I summarizes all system parameters in our experiments. Unless otherwise stated, our results assume a  $32nm$  HP circuit technology, a  $10.2ns$  wake-up mode, and an out-of-order CMP. The following overheads were considered when calculating the reported results: core wake-up energy, core wake-up delay, core pipeline-refill latency, retention overhead of live-slave retention cells, SRAM leakage during source biasing mode of operation, and voltage noise safety.

We simulate the system with the GEM5 simulator [5]. GEM5 features cycle-level models of an out-of-order core, the cache hierarchy, and the interconnect. We integrate GEM5 with DRAMSim2 [1] to provide cycle-level modeling of the memory subsystem including the memory controller, DRAM modules, and shared channels used for communication. We

parameter	value	notes
EV6 core model	DEC-Alpha EV6	
EV6 core clock	3.3GHz-1.9GHz	
EV6 execution	6-way out-of-order	
EV6 functional units	6ALU,2IMULT 2FPALU	
ICache/DCache	32KB-8way 1 cycle	
L2 Cache	256KB-8way 4ns	Private per core
L3 Cache	8MB-16way 13ns	Shared
Memory	DDR3 2GB 50ns	
Core-to-L1 token latency	0.5ns	controller delays
Core-to-L2 token latency	4.5ns	controller delays
Core-to-L3 token latency	17.5ns	controller delays
Core-to-centralized-WUC latency	5ns	controller delays
PPGS wake-up modes	4.5ns-16.9ns	SPICE
EV6 pipeline refill latency	2.12ns	7 pipeline stages
EV6 core wake-up energy (EWE)	15,358pJ	Charge cells
EV6 leakage power (ELP)	0.916 Watts	McPAT [17]
EV6 PG leakage reduction (ELPR)	97.65%	[7]
EV6 PG break even point	17.17ns	$EWE/(ELPR*ELP)$
EV6 DFLT core wake-up latency	10.2ns	SPICE
EV6 FUPG wake-up energy	9641pJ	McPAT, ITRS [2]
EV6 FUPG wake-up latency	6.4ns	SPICE

TABLE I  
SYSTEM CONFIGURATION VALUES

use McPAT [17] to model dynamic power, leakage power, peak power, and area. We update McPAT's *technology.cc* file to accurately reflect the ITRS 2010 update report [2].

#### A. Adapting to Memory Contention.

We show how TAP adapts to a system facing increased memory contention for the multi-threaded memory benchmark *STREAM* running on a CMP with up to 32 cores. *STREAM* is a memory-intense benchmark used to measure sustained memory bandwidth and computation rates for simple vector kernels [3]. We modify *STREAM* to act as a parallel memory benchmark such that more threads cause more simultaneous requests to the memory subsystem. Increasing *STREAM*'s thread count causes more queuing of memory requests, longer delay per request and more frequent stalling of each core. A good power-gating technique should be able to power gate the core more often to reduce power consumption.

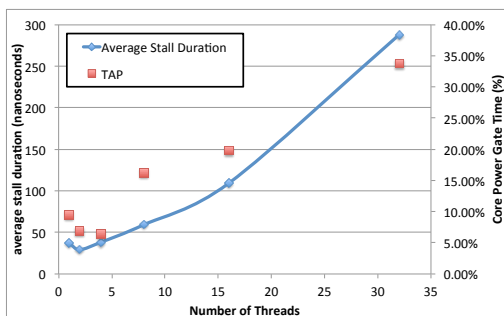


Fig. 4. TAP adapting to increasing memory contention. The left y-axis shows the duration of the stall in nanoseconds, while the right y-axis shows the percentage of time that TAP can power gate as a function to the number of *STREAM* threads.

Figure 4 depicts both average duration of core stalls and the percentage of total simulation time TAP power gates the core. First, we note that as the number of threads increases, the average duration of a core stall increases. From 1 to 32 threads, average core-stall duration increases<sup>2</sup> by 7.82 $\times$  from

<sup>2</sup>From 1 to 2 threads, core-stall duration decreases. This happens because more threads increase the amount of available cache (more cores) while increasing the number of simultaneous memory requests.

36.77ns to 287.63ns. The increase in the average core-stall duration is caused by more threads making parallel requests to the memory subsystem at once, fewer row-buffer hits, and memory channel contention. As cores experience increased memory latency, TAP power gates the core longer. From 1 to 32 threads, TAP power gates cores 3.69 $\times$  longer from 9.08% to 33.50% of simulation time. However, TAP power gates its core less for 4 threads than for 2 even though average core-stall time increases. This is because TAP uses a conservative lower-bound estimate of memory-response time and does not account for all memory scheduling possibilities. The benefit of this conservative policy is the absence of any application performance hit.

#### B. Distributed, Staggered Wake Up

For a 16-core system with multiple cores waking up simultaneously, voltage noise on the power distribution network can cause unsafe voltage drops on neighboring active cores. By introducing a sub-nanosecond stagger between any two adjacent cores waking up, worst-case inrush current and resulting voltage noise are reduced. The result is a faster wake-up mode and increased energy savings on the chip. Hence cores should wake up staggered by at least a fraction of a nanosecond.

Figure 5 shows SPICE simulation of the effect of stagger on core wake-up latency for no-stagger, 0.3ns-stagger, 0.6ns-stagger, and 0.9ns-stagger for CMPs composed of 16 EV6 cores as 0 to 14 cores are idle. Staggered wake up can reduce the variance between the min and max wake-up latency when most cores are actively executing or waking-up. For the 16-core CMP with no cores idle and no stagger, the max and min wake-up latency is 18.4ns and 9.7ns, whereas, a staggered wake up of 0.9ns decreases the max and min values to 10.7ns and 8.6ns respectively. As more cores go idle, staggered wake up has less impact on reducing the variance of wake-up latency because cores are less likely to interfere with each other, and the core's location becomes the dominant factor in wake-up latency. For example, when 14 cores are idle in a 16-core CMP, the min and max wake-up latencies are both 4.5ns and 9.1ns, respectively, in both the no-stagger and 0.9ns-stagger cases. Lastly, stagger reduces the maximum wake-up latency as more cores are active. An example of this is for the 16-core case when no core is idle; maximum wake-up latency is 18.4ns without stagger and 10.7ns with 0.9ns stagger, a reduction of 58.2%. Thus, stagger relaxes the guardband on wake-up latency.

Figure 6 shows the energy impact of staggered wake up for a CMP with 16 EV6 cores. The largest improvement in energy savings is 3.14% for *mcf* as energy savings increase from 18.92% to 22.06% for staggered wake ups of 0.0ns and 0.9ns respectively. On average, energy savings go from 3.52% to 4.34% as stagger increases from 0ns to 0.9ns. Staggered wake up does not cause any decrease in energy savings. Although cores may have to wake up at slightly different times, the savings in latency with which they wake up is much greater than the stagger offset of 0.9ns.

To have cores safely and reliably use staggered wake up, a controller or scheme is required to control the wake-up behavior of all cores. A centralized WUC would be able to use more aggressive wake-up modes and power gate the core for longer, but such a design does not scale to many cores.



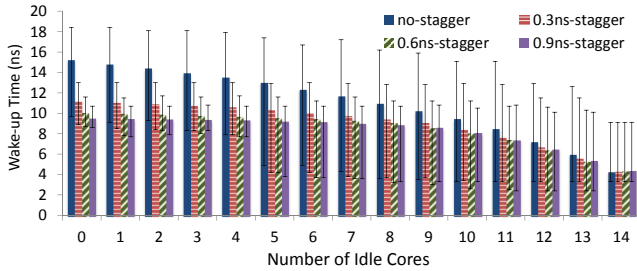


Fig. 5. The improvement in core wake-up latency with increased stagger for a 16 EV6 core CMP as 0 to 14 cores are idle. The average latency is shown with bars denoting the min and max safe wake-up modes.

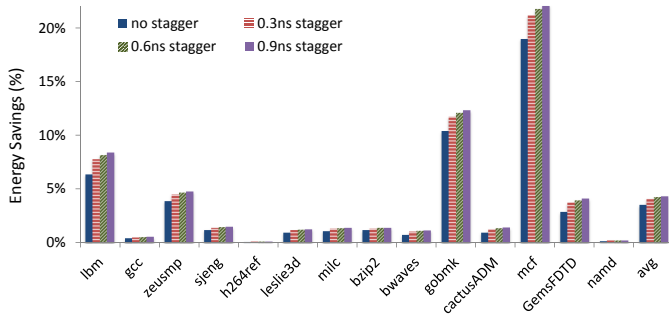


Fig. 6. The improvement in energy savings with staggered wake ups in an 16 EV6 core CMP. Benchmarks with less than 0.2% energy savings filtered.

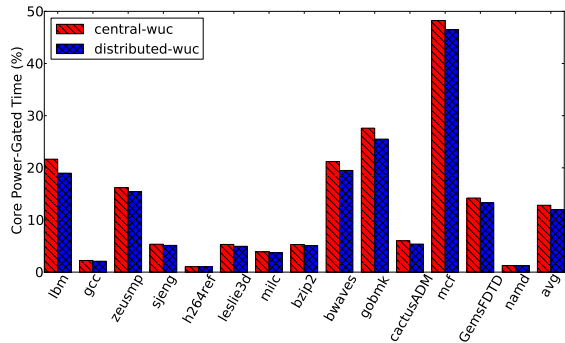


Fig. 7. Comparison of power-gated time using staggered wake up with a centralized WUC versus a distributed WUC.

By contrast, the distributed wake-up control of Section III can scale to many-core designs, but sacrifices some power-gating time waiting for an assigned wake-up slot. Figure 7 shows the difference of application power-gated time for a subset of the SPEC2006 benchmarks<sup>3</sup> in a 16-core architecture. Our simulations show that the distributed wake-up control has a maximum decrease in power-gated time of 2.68% (0.93% on average) for the benchmark *lbn*. This indicates that the difference in energy savings between a distributed scheme and a centralized WUC would be negligible, but that a distributed scheme scales to many-core architectures.

## V. CONCLUSIONS

With each new generation of microprocessors, leakage power becomes an increasingly dominant issue. In this paper, we extend TAP [11] to many-core designs through a distributed, staggered wake-up scheme. Our system can adapt to different levels of memory contention by increasing power-gated time by  $3.69\times$  of total execution time as the number of

threads increases from 1 to 32. We demonstrate that a wake-up stagger of  $0.9ns$  reduces core wake-up latency by up to 58.2% (7.7ns). Last, a distributed wake-up scheme that uses staggered wake-up can power gate cores for 99.07% of the time achieved by an ideal centralized scheme, while still being able to scale to many-core designs.

## VI. ACKNOWLEDGMENTS

This research was supported by the MARCO FCRP (GSRC and MuSyC centers), Qualcomm, Oracle, and NSF grants SHF-0916127, SHF-1218666, SHF-1116667 and CCF-1162085. We thank Kwangok Jeong for his early work on this project [10].

## REFERENCES

- [1] *DRAMSim2*, 2011, <http://www.ece.umd.edu/dramsim/>.
- [2] *International Technology Roadmap for Semiconductors*, 2010, <http://www.itrs.net/>.
- [3] *STREAM: Sustainable Memory Bandwidth in High Performance Computers*, 2011, <http://www.cs.virginia.edu/stream/>.
- [4] K. Agarwal, H. Deogun, D. Sylvester and K. Nowka, "Power Gating with Multiple Sleep Modes", *Proc. ISQED*, 2006, pp. 633–637.
- [5] N. L. Binkert, R. G. Dreslinski, L. R. Hsu, K. T. Lim, A. G. Saidi and S. K. Reinhardt, "The M5 Simulator: Modeling Networked Systems", *IEEE Micro* 26(4) (2006), pp. 52–60.
- [6] M. H. Chowdhury, J. Gjanci and P. Khaled, "Innovative Power Gating for Leakage Reduction", *Proc. ISCAS*, 2008, pp. 1568–1571.
- [7] D. Flynn, R. Aitken, A. Gibbons and K. Shi, *Low Power Methodology Manual*, Springer, 2007.
- [8] M. Horiguchi, T. Sakata and K. Itoh, "Switched-Source-Impedance CMOS Circuit for Low Standby Subthreshold Current Giga-Scale LSI's", *IEEE JSSC* 28(11) (1993), pp. 1131–1135.
- [9] Z. Hu, A. Buyuktosunoglu, V. Srinivasan, V. Zyuban, H. Jacobson and P. Bose, "Microarchitectural Techniques for Power Gating of Execution Units", *Proc. ISLPED*, 2004, pp. 32–37.
- [10] K. Jeong, A. B. Kahng, S. Kang, T. S. Rosing and R. Strong, "MAPG: Memory Access Power Gating", *Proc. DATE*, 2012, pp. 1054–1059.
- [11] A. B. Kahng, S. Kang, T. S. Rosing and R. Strong, "TAP - Token-Based Adaptive Power Gating", *Proc. ISLPED*, 2012, pp. 203–208.
- [12] S. Kim, S. V. Kosonocky and D. R. Knebel, "Understanding and Minimizing Ground Bounce during Mode Transition of Power Gating Structures", *Proc. ISLPED*, 2003, pp. 22–25.
- [13] R. Kumar and G. Hinton, "A Family of 45nm IA Processors", *Proc. ISSCC*, 2009, pp. 58–59.
- [14] H. Li, C. Cher, T. N. Vijaykumar and K. Roy, "VSV: L2-miss-driven variable supply-voltage scaling for low power", *Proc. of the IEEE International Symposium on Microarchitecture*, 2003, pp. 19–28.
- [15] Y. Lin, C. Yang, J. Huang and N. Chang, "Memory Access Aware Power Gating for MPSoCs", *Proc. ASP-DAC*, 2012, pp. 121–126.
- [16] J. Leverich, M. Monchiero, V. Talwar, P. Ranganathan and C. Kozyrakis, "Power Management of Datacenter Workloads Using Per-Core Power Gating", *IEEE Computer Architecture Letters* 8(2) (2009), pp. 48–51.
- [17] S. Li, J. H. Ahn, R. D. Strong, J. B. Brockman, D. M. Tullsen and N. P. Jouppi, "McPAT: An Integrated Power, Area, and Timing Modeling Framework for Multicore and Manycore Architectures", *Proc. MICRO*, 2009, pp. 469–480.
- [18] A. Lungu, P. Bose, A. Buyuktosunoglu and D. J. Sorin, "Dynamic Power Gating with Quality Guarantees", *Proc. ISLPED*, 2009, pp. 377–382.
- [19] N. Madan, A. Buyuktosunoglu, P. Bose and M. Annavaram, "A Guarded Power Gating for Multi-Core Processors", *Proc. HPCA*, 2011, pp. 291–300.
- [20] E. Perelman, G. Hamerly, M. Van Biesbrouck, T. Sherwood and B. Calder, "Using SimPoint for Accurate and Efficient Simulation", *Proc. SIGMETRICS*, 2003, pp. 318–319.
- [21] H. Qin, Y. Cao, D. Markovic, A. Vladimirescu and J. Rabaey, "SRAM Leakage Suppression by Minimizing Standby Supply Voltage", *Proc. ISQED*, 2004, pp. 55–60.
- [22] A. Rogers, D. Kaplan, E. Quinnell and B. Kwan, "The Core-C6 (CC6) Sleep State of the AMD Bobcat x86 Microprocessor", *Proc. ISLPED*, 2012, pp. 367–372.
- [23] Y. Shin, J. Seomun, K.-M. Choi and T. Sakurai, "Power Gating: Circuits, Design Methodologies, and Best Practice for Standard-Cell VLSI Designs", *ACM TODAES* 15(4) (2010), pp. 1–37.
- [24] H. Singh, K. Agarwal, D. Sylvester and K. Nowka, "Enhanced Leakage Reduction Techniques Using Intermediate Strength Power Gating", *IEEE TVLSI* 15(11) (2007), pp.1215–1224.
- [25] O. Wechsler, "Setting New Standards for Energy-Efficient Performance", *Technology@Intel Magazine*, 2006.
- [26] Z. Zhang, X. Kavousianos, K. Chakrabarty and Y. Tsiatouhas, "A Robust and Reconfigurable Multi-mode Power Gating Architecture", *Proc. VLSID*, 2011, pp. 280–285.
- [27] H. Zheng and Z. Zhu, "Power and Performance Trade-Offs in Contemporary DRAM System Designs for Multicore Processors", *IEEE TC* 59(8) (2010), pp. 1033–1046.

<sup>3</sup>We use the Simpoint methodology [20] with 100M-instruction representative regions for each application.