

# Exploiting Spatial Locality for Energy-Optimized Compressed Caching

Somayeh Sardashti

PhD Advisor: Professor David A. Wood

University of Wisconsin-Madison

## Abstract

Last-level caches (LLCs) play a crucial role in reducing multicore system energy by filtering out expensive accesses to main memory. Cache compression can increase effective LLC capacity and reduce misses. However, previous designs limit compression benefits caused by internal fragmentation, limited tags, and energy-expensive recompaction when a block's size changes. In this work, we propose decoupled compressed cache (DCC), which uses decoupled superblocks and noncontiguous sub-block allocation to decrease tag overhead without increasing internal fragmentation and to eliminate the need for energy-expensive recompaction. We also demonstrate a practical design based on a recent commercial LLC design.

## 1. Problem and Motivation

Future computer systems face continuing power and energy challenges as the power per transistor scales more slowly than transistor density [14]. Caches, especially last-level caches (LLCs), have long been used to reduce effective memory latency and increase effective bandwidth. They also play an increasingly important role in reducing memory system energy. Keckler [15] shows that last-level caches (LLCs) are especially important, since obtaining operands of a double-precision multiply-add from off-chip memory requires approximately 200x the energy of the operation. Thus improving effective cache utilization is important not only for system performance, but also for system energy.

Increasing LLC size can improve performance for most workloads, but the improvement comes at significant area cost. For example, the well-known “square root” power law [15] predicts that doubling LLC size will reduce misses by ~30%, on average. But it obviously doubles LLC area, which already accounts for 15–30% of the die area of most processors [8].

Cache compression, however, seeks to increase effective cache size—by compressing and compacting cache blocks—while incurring small area overheads [2][17][18][4]. For example, previously proposed techniques have the potential to double effective LLC capacity, while increasing LLC area by only ~8%.

Although some data (and most instructions) are

difficult to compress, most workloads are highly compressible. In this article, we use C-PACK+Z, a dictionary-based algorithm<sup>4</sup> with nine-cycle decompression latency. C-PACK+Z achieves an average compression ratio (that is, the original size over compressed size) of 3.9. Thus, compression has the potential to nearly quadruple cache size (shown as “Ideal” in Figure 1).

Despite compression potentials, previous compressed cache designs fail to achieve this potential for three main reasons. First, caches must compact compressed blocks into sets, which introduces an internal fragmentation problem. In Figure 1, BytePack represents an idealized compressed cache with infinite tags, which compacts compressed blocks on arbitrary byte boundaries. BytePack degrades normalized effective capacity to 3.1 on average. Second, practical compressed caches introduce another internal fragmentation problem by compacting compressed blocks into one or more sub-blocks, rather than storing compressed data on arbitrary byte boundaries [2]. Variable-size compression (VSC) techniques relax the mapping constraint between tags and data and compact compressed blocks into a variable number of contiguous sub-blocks [2]. The column labeled VSC-Inf in Figure 1 illustrates that compacting compressed blocks into zero to four 16-byte sub-blocks (but with infinite tags per set) degrades normalized effective capacity from 3.1 to 2.6 on average. Third, practical compressed caches have a fixed number of tags per set. The remaining columns in Figure 1 illustrate that reducing the number of tags, from infinite to a more practical two times the baseline, degrades the average normalized effective capacity from 2.6 to 1.7. Furthermore, VSC is not energy efficient. It must repack

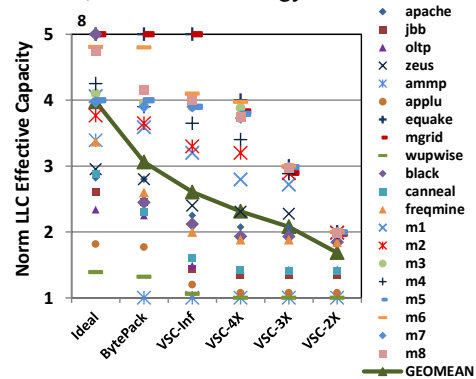


Figure 1. Normalized effective capacity of different compressed cache designs

the sub-blocks in a set whenever a block’s size changes to make contiguous free space. This action can increase LLC dynamic energy by a factor of nearly three, on average.

In this work, we propose Decoupled Compressed Cache that exploits spatial locality to address the limits with previous compressed caches. It improves both the performance and energy-efficiency of cache compression.

## 2. Background and Related Work

**Exploiting Spatial Locality in Caches:** This work builds upon previous dual-grain caches namely the Region Tracker [19], the sectored cache [20], sector pool cache [21], and the decoupled sectored cache [3]. RegionTracker also manages cache at dual-granularities of memory regions (e.g., 1KB) and cache blocks [19]. Unlike our proposal, Region Tracker [19] aims facilitating collection of coarse-grain information.

**Cache Compression:** Hallnor et al. extend their earlier indirect index cache [23] to support compression (IIC-C) [22]. IIC-C uses a software-managed hash table to provide full associativity and forward pointers to associate tags with variable number of sub-blocks anywhere in the data array [22], eliminating the need for repacking. However, for an 8MB LLC with 64-byte blocks, 16-byte sub-blocks, and doubled number of tags, their scheme incurs about 24% area overhead (26% considering (de-)compressors), while it at most doubles effective capacity. Further increasing the number of tags will make its area overhead even worse. J. S. Lee, et al. [24] compress block pairs and store them in a single line if both lines compress by 50% or more. In this way, they free a cache block in an adjacent set; however, they need to check two sets for a potential hit on every access, which increases power overheads. In addition, their technique limits compressibility by failing to take advantage of lines that compress by less than 50%.

Compression is also used at main memory [25][26]. IBM’s MXT effectively doubles main memory capacity [25]. The decoupled zero-compressed memory [27] manages the main memory as a decoupled sectored set-associative cache. It detects null blocks to improve performance, which is more limited than compression based schemes.

## 3. Approach and Uniqueness

We propose decoupled compressed cache (DCC), a technique that uses decoupled superblocks (also known as sectors<sup>3</sup>) to increase the maximum effective capacity to four times the uncompressed capacity, while using area

overhead comparable to previous cache-compression techniques. DCC uses superblocks—four aligned, contiguous cache blocks that share a single address tag—to reduce tag overhead. Each 64-byte block in a superblock is compressed separately and then compacted into zero to four 16-byte sub-blocks (or segments). DCC decouples the address tags by increasing the number of tags and allowing any sub-block in a set to map to any tag in that set to reduce fragmentation within a superblock [3].

Decoupling allows sub-blocks of a block to be noncontiguous, thereby eliminating the recompact overheads of previous variable-size compressed caches [2]. An optimized co-compacted DCC (Co-DCC) design further reduces internal fragmentation (and increases effective capacity) by allocating the compressed blocks from a superblock into the same set of data sub-blocks.

This article makes the following contributions:

- DCC uses decoupled superblocks to increase the effective number of tags with low overhead.
- DCC stores compressed data in noncontiguous sub-blocks to eliminate recompact overheads when a block’s compressed size changes.
- DCC provides more effective capacity, on average, than a conventional cache of twice the size, while slightly increasing cache area. Viewed another way, DCC allows a designer to get approximately the same cache performance with about half the area.
- Co-DCC further reduces internal fragmentation by compacting the blocks of a superblock and allocating them into the same set of data sub-blocks.
- We also present a concrete design for Co-DCC and show how it can be integrated into a recent commercial LLC design with little additional complexity.

### 3.1 Decoupled compressed cache

DCC uses decoupled superblock tags to improve cache compression in two ways. First, superblocks reduce tag overhead, permitting more tags per set for comparable overhead. Second, decoupling tags and data reduces internal fragmentation and, importantly, eliminates recompact when the size of a compressed block changes.

Figure 2 shows how DCC exploits superblocks and manages the cache at three granularities: coarse-grained superblocks, singular cache blocks, and fine-grained sub-blocks. DCC tracks superblocks, which are groups of

aligned, contiguous cache blocks (Figure 2d), while it compresses and stores each cache block as a variable number of sub-blocks. Figure 2a shows the key components of DCC for a small, two-way set associative cache with four-block superblocks, 64-byte blocks, and 16-byte sub-blocks. DCC consists of a tag array, a sub-blocked back pointer array, and a sub-blocked data array. DCC is indexed using the superblock address bits (“Set Index” in Figure 2e).

DCC tracks superblocks to fit more compressed blocks into the cache while limiting tag area overhead. DCC explicitly tracks superblocks using a largely conventional superblock tag array. Each tag entry (Figure 2b) consists of one tag per superblock and per-block coherence (C-state) and compression (Comp) states. Because blocks of a superblock share an address tag, the tag array can map more blocks compared with the same size conventional cache without incurring high area overhead. DCC holds as many superblock tags as the maximum number of uncompressed blocks that can be stored. For example, Figure 2a shows a two-way associative cache with four-block superblocks. Each set in the tag array can map eight blocks (that is, 2 superblocks  $\times$  4 blocks/superblock), while a maximum of two uncompressed blocks can fit in each set. In the worst-case scenario—when there is no spatial locality (that is, all singletons) or when cached data is uncompressible—DCC can still utilize all the cache data space, for example, by tracking two singletons per set.

Although DCC tracks blocks of a superblock using one tag entry, it allocates or evicts the blocks to and from the data array separately. The data array is a mostly conventional cache data array, organized into sub-blocks. DCC compacts compressed blocks into a variable number of noncontiguous sub-blocks in the sub-blocked data array. Figure 2a shows block A0 compressed into two sub-blocks (A0.1 and A0.0), which are stored in sub-blocks 5 and 1 in the data array. DCC decouples sub-

blocks from the superblock tag using a back pointer array as a level of indirection. Each back pointer entry corresponds to one sub-block in the data array and identifies the owner block (Figure 2c). The back pointer array slightly increases LLC area; however, it enables low-overhead, variable-size compression. DCC’s decoupled design allows a block’s sub-blocks to be noncontiguous, thus eliminating the need for recompaction when a block’s size changes.

Co-DCC optimizes DCC to further reduce internal fragmentation. Co-DCC treats blocks from the same superblock as one large block and dynamically allocates them into the same set of data sub-blocks, thereby reducing internal fragmentation within sub-block boundaries. Co-DCC increases overhead and complexity in exchange for better cache compression.

Figure 2f illustrates the DCC lookup procedure. On a cache lookup, the tag array and the back pointer array are accessed in parallel. In the common case of a cache hit, both the block and its corresponding superblock are found available (that is, tag matched and block is valid). In the event of a cache hit, the result of the tag array and the back pointer array lookup determines which sub-blocks of the data array belong to the accessing block. On a read, the corresponding sub-blocks are then read out of the data array and decompressed. On a write, the new, compressed size might be larger, resulting in a block (or a superblock) eviction if sufficient space is not available. On the other hand, in case of a cache miss, DCC allocates the compressed block in the data array. If its superblock is not available, DCC allocates it first in the tag array.

### 3.2 A practical design for DCC

DCC can be integrated into the LLC of a recent commercial design with relatively little additional complexity and, more importantly, no need for an

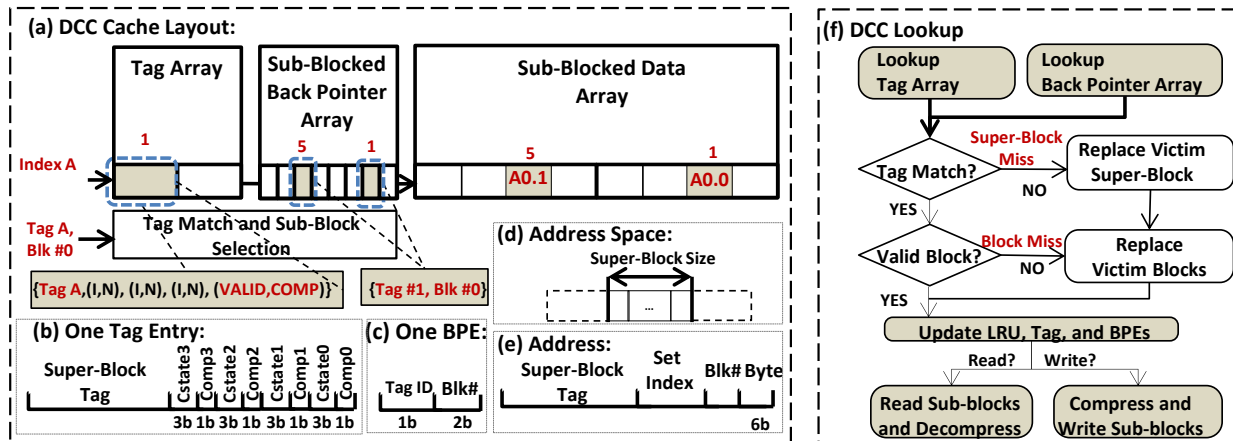


Figure 2. A decoupled compressed cache. DCC cache layout (a); one tag entry (b); one back pointer entry (BPE) (c); address space (d); address (e); and DCC lookup process (f).

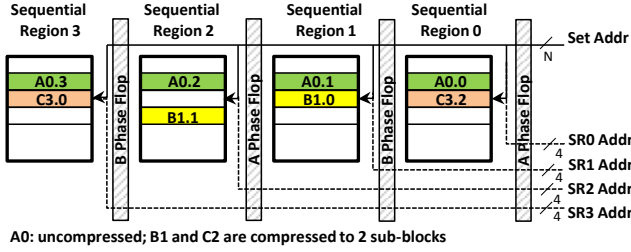


Figure 3. DCC data array organization.

alignment network. The AMD Bulldozer processor implements an 8-Mbyte LLC that is broken into four 2-Mbyte subcaches. Each subcache consists of four banks that can independently service cache accesses[5]. Figure 3 illustrates the data array of one bank in LLC and shows how it is divided into four sequential regions (SRs). Each sequential region runs one phase (that is, half a cycle) behind the previous region and contains a quarter of a cache block (that is, 16 bytes). Figure 3 shows how block A0’s four 16-byte sub-blocks (A0.0 to A0.3) are distributed to the same row in each sequential region. Each subsequent sequential region receives the address half a cycle later and takes half a cycle longer to return the data. Thus, a 64-byte block is returned in a burst of four cycles on the same data bus. For example, A0.1 is returned one cycle after A0.0.

DCC requires only a small change to the data array to allow noncontiguous sub-blocks. In Figure 3, block B1 is compressed into two sub-blocks (B1.0 and B1.1), which are stored in sequential regions 1 and 2, but not in the same row. To select the correct sub-block, DCC must send additional address lines (4 bits for a 16-way associative cache) to each sequential region (illustrated by the dotted lines in Figure 3). DCC must also enforce the constraint that a compressed block’s sub-blocks are allocated to different sequential regions to prevent sequential region conflicts.

Figure 4b illustrates DCC timing when reading block B1. The back pointer array is accessed in parallel with the tag array. The sub-block selection logic then finds the back pointer entries corresponding to this block using its block ID (derived from its address) and the matched tag ID, which is found by the tag match logic. The sub-block selection logic can only be partially overlapped with the tag match logic because it needs the matched tag ID. To calculate the latency overhead of the sub-block selection, we implemented the tag match and the subselection logic in Verilog, synthesized in 45 nm, and scaled to 32 nm [6]. The sub-block selection logic adds less than half a cycle to the critical path, which we conservatively assume increases the access latency by one cycle.

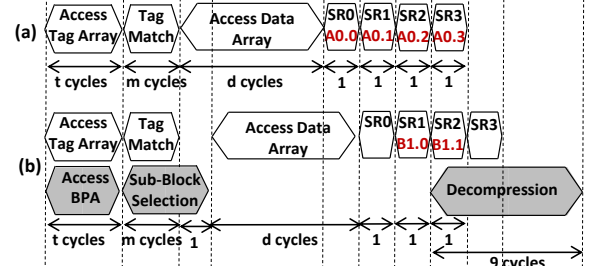


Figure 4. Timing of a conventional cache (a) and DCC (b)

## 4. Results and Contributions

We evaluate DCC using a full-system simulator based on GEMS [7]. We model a multicore system with eight out-of-order cores; per-core private 32-Kbyte, 8-way L1 instruction and data caches; per-core private 256-Kbyte, 8-way L2 caches; and one shared 8-Mbyte, 16-way L3 cache [8]. We use CACTI [9] to model power at 32 nm. We also use a detailed DRAM power model based on Micron Technology’s power model [10]. In this section, we report total system energy, which includes energy consumption of processors (cores and caches), on-chip network, and off-chip memory. For DCC and Co-DCC, we use four-block superblocks, 64-byte blocks, and 16-byte sub-blocks. With these parameters, DCC has similar area overhead as FixedC (~8% LLC area overhead), which doubles the number of tags and compresses a block to half, if possible, and VSC-2X, which doubles tags but compresses a block into zero to four 16-byte sub-blocks.

Our evaluations use representative multithreaded and multiprogrammed workloads from commercial workloads (*apache*, *job*, *oltp*, *zeus*) [11], SPEC-OMP (*ammp*, *applu*, *equake*, *mgrid*, *wupwise*) [12], Parsec (*blackscholes*, *cannal*, *freqmine*) [13], and mixes of SPEC CPU2006 benchmarks denoted as m1 to m8 (*bzip2*, *libquantum-bzip2*, *libquantum*, *gcc*, *astar-bwaves*, *cactus-mcf-milc-bwaves*, *gcc-omnetpp-mcf-bwaves-lbm-milc-cactus-bzip*, *omnetpp-lbm*).

### 4.1 Improved cache efficiency

Compressed caches improve the cache’s effective capacity by fitting more blocks into the same space. They can achieve the benefits of larger cache sizes with lower area and power overheads.

**Result 1:** By exploiting spatial locality, DCC achieves on average 2.2× (and up to 4×) higher LLC effective capacity compared to the baseline, resulting in 18 percent lower LLC miss rate on average and up to 38 percent lower LLC miss rate.

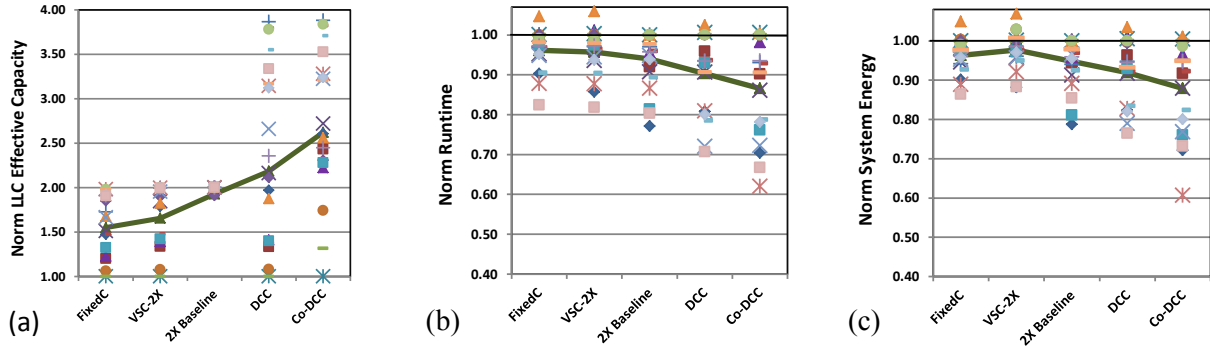


Figure 5. Normalized LLC effective capacity (a); normalized runtime (b); normalized total system energy (c).

**Result 2:** Co-DCC further improves the effective cache capacity by reducing internal fragmentation within data sets. It achieves on average 2.6 $\times$  (and up to 4 $\times$ ) higher effective capacity and 24 percent (up to 42 percent) lower LLC miss rate.

**Result 3:** DCC and Co-DCC provide significantly higher effective cache capacity and lower miss rate than FixedC and VSC-2X. DCC and Co-DCC also perform better on average than a cache at twice the capacity (2 $\times$  baseline) with much lower area overhead.

Figure 5a shows LLC effective capacity of different techniques normalized to baseline. We calculate the effective cache capacity by counting valid LLC cache blocks periodically. DCC can significantly improve LLC effective capacity and LLC miss rate (misses per kilo executed instructions [MPKI]) for many applications by fitting more compressed blocks. DCC benefits differ per workload, depending on workload sensitivity to cache capacity, compression ratio, and spatial locality. It achieves the greatest benefit for cache-sensitive workloads with good compressibility and spatial locality (such as *apache* and *omnetpp-lbm/m8*). Workloads with low spatial locality (such as *canneal*) or low compression ratio (such as *wupwise*) observe lower improvements. Cache-insensitive workloads (such as *blackscholes*) also do not benefit from compression.

## 4.2 Overall performance and energy

By improving LLC utilization and reducing accesses to the main memory (that is, the lower LLC miss rate), DCC and Co-DCC significantly improve system performance and energy.

**Result 4:** DCC and Co-DCC improve LLC efficiency and boost system performance by 10 percent (up to 29 percent) and 14 percent (up to 38 percent) on average, respectively.

**Result 5:** DCC and Co-DCC save on average 8 percent

(up to 24 percent) and 12 percent (up to 39 percent) of system energy, respectively, because of shorter runtime and fewer accesses to the main memory.

**Result 6:** DCC and Co-DCC achieve 2.5 $\times$  and 3.5 $\times$  higher performance improvements, respectively, and 2.2 $\times$  and 3.3 $\times$  higher system energy improvements compared with FixedC and VSC-2X.

**Result 7:** DCC and Co-DCC also improve LLC dynamic energy by about 50 percent on average by accessing fewer bytes. However, VSC-2X hurts LLC dynamic energy for the majority of our workloads because of its need for energy-expensive recombinations.

Figure 5b shows that DCC outperforms baseline, FixedC, VSC-2X, and 2 $\times$  baseline by effectively more than doubling the cache capacity. DCC and Co-DCC also improve system energy owing to shorter runtime and fewer accesses to the main memory. Figure 5c shows the total system energy of different techniques. DCC and Co-DCC significantly reduce the main memory dynamic energy by reducing the number of cache misses, which contributes to greater system energy improvements as well. Unlike VSC-2X, which hurts LLC dynamic energy because of recombination, DCC and Co-DCC eliminate the need for recombination and can even save LLC dynamic energy by accessing fewer bytes when reading or writing compressed data.

## 5. Conclusions and Future Work

DCC demonstrates the potential for compression to increase the effective capacity of last-level caches, improving both performance and energy efficiency. Alternatively, DCC can reduce the chip area required for a given cache capacity, thereby reducing implementation cost. Future work should explore algorithms that better compress instructions and floating point data, as well as extending compression to all levels of the cache and memory hierarchy.

## References

- [1] S. Sardashti and D. Wood. "Decoupled Compressed Cache: Exploiting Spatial Locality for Energy-Optimized Compressed Caching," *Proc. 46th Ann. IEEE/ACM Int'l Symp. Microarchitecture*, 2013.
- [2] A. Alameldeen and D. Wood. "Adaptive Cache Compression for High-Performance Processors," *Proc. 31st Ann. Int'l Symp. Computer Architecture*, 2004.
- [3] A. Sez nec, "Decoupled Sector ed Caches: Conciliating Low Tag Implementation Cost and Low Miss Ratio," *Proc. 21st Ann. Int'l Symp. Computer Architecture*, 1994.
- [4] X. Chen et al., "C-Pack: A High-Performance Microprocessor Cache Compression Algorithm," *IEEE Trans. VLSI Systems*, vol. 18, no. 18, 2010.
- [5] D. Weiss et al., "An 8MB Level-3 Cache in 32nm SOI with Column-Select Aliasing," *Proc. Solid-State Circuits Conf.*, 2011.
- [6] *International Technology Roadmap for Semiconductors, 2010 Update*, ITRS, 2011; [www.itrs.net](http://www.itrs.net).
- [7] M. Martin et al., "Multifacet's General Execution-Driven Multiprocessor Simulator (GEMS) Toolset," *Computer Architecture News*, 2005, vol. 33, no.4.
- [8] "4th Generation Intel Core i7 Processors," Intel Corporation; [www.intel.com/products/processor/corei7](http://www.intel.com/products/processor/corei7).
- [9] "CACTI: An Integrated Cache and Memory Access Time, Cycle Time, Area, Leakage, and Dynamic Power Model," HP Labs Research; [www.hpl.hp.com/research/cacti](http://www.hpl.hp.com/research/cacti).
- [10] "Calculating Memory System Power for DDR3," tech. note TN-41-01, Micron Technology, 2007.
- [11] A. Alameldeen et al., "Simulating a \$2M Commercial Server on a \$2K PC," *IEEE Computer*, vol. 36, no. 2, 2003.
- [12] V. Aslot et al., "SPECComp: A New Benchmark Suite for Measuring Parallel Computer Performance," *Proc. Int'l Workshop OpenMP Applications and Tools: OpenMP Shared Memory Parallel Programming*, 2001, pp. 1-10.
- [13] C. Bienia and K. Li, "PARSEC 2.0: A New Benchmark Suite for Chip-Multiprocessors," *Proc. 5th Ann. Workshop Modeling, Benchmarking and Simulation*, 2009, pp. 47-55.
- [14] Dennard R. et al. 1974. Design of Ion-Implanted MOSFET's with Very Small Physical Dimensions. *IEEE Journal of Solid-State Circuits*.
- [15] Keckler, S. 2011. Life After Dennard and How I Learned to Love the Picojoule. In Proceedings of the 44th Annual IEEE/ACM International Symposium on Microarchitecture.
- [16] Hartstein, A. et al. 2008. On the Nature of Cache Miss Behavior: Is It  $\sqrt{2}$ ? *J. Instruction-Level Parallelism* 10.
- [17] Arelakis, A., Stenström, P. 2012. A Case for a Value-Aware Cache. *IEEE Computer Architecture Letters*.
- [18] Kim, N. et al. 2002. Low-Energy Data Cache Using Sign Compression and Cache Line Bisection. Second Annual workshop on Memory Performance Issues.
- [19] Jason Zebchuk , Elham Safi , Andreas Moshovos, A Framework for Coarse-Grain Optimizations in the On-Chip Memory Hierarchy, Proceedings of the 40th Annual IEEE/ACM International Symposium on Microarchitecture, p.314-327, December 01-05, 2007
- [20] J. S. Liptay, Structural aspects of the system/360 model 85: II the cache, *IBM Systems Journal*, v.7 n.1, p.15-21, 1968.
- [21] Jeffrey B. Rothman , Alan Jay Smith, The pool of subsectors cache design, Proceedings of the 13th international conference on Supercomputing, , 1999.
- [22] Erik G. Hallnor , Steven K. Reinhardt, A Unified Compressed Memory Hierarchy, Proceedings of the 11th International Symposium on High-Performance Computer Architecture, p.201-212, February 12-16, 2005.
- [23] Erik G. Hallnor , Steven K. Reinhardt, A fully associative software-managed cache design, Proceedings of the 27th annual international symposium on Computer architecture, p.107-116, June 2000.
- [24] Jang-Soo Lee , Won-Kee Hong , Shin-Dug Kim, An on-chip cache compression technique to reduce decompression overhead and design complexity, *Journal of Systems Architecture: the EUROMICRO Journal*, 2000.
- [25] Bulent Abali , Hubertus Franke , Xiaowei Shen , Dan E. Poff , T. Basil Smith, Performance of Hardware Compressed Main Memory, Proceedings of the 7th International Symposium on High-Performance Computer Architecture, p.73, January 20-24, 2001.
- [26] Magnus Ekman , Per Stenstrom, A Robust Main-Memory Compression Scheme, *ACM SIGARCH Computer Architecture News*, v.33 n.2, p.74-85, May 2005.
- [27] Julien Dusser , André Sez nec, Decoupled zero-compressed memory, Proceedings of the 6th International Conference on High Performance and Embedded Architectures and Compilers, January 24-26, 2011.