

Skewed Caches from a Low-Power Perspective

Mathias Spjuth, Martin Karlsson and Erik Hagersten
Uppsala University
Information Technology
Department of Computer Systems
P.O. Box 325, SE-751 05 Uppsala, Sweden
martink@it.uu.se

ABSTRACT

The common approach to reduce cache conflicts is to increase the associativity. From a dynamic power perspective this associativity comes at a high cost. In this paper we present miss ratio performance and a dynamic power comparison for set-associative caches, a skewed cache and also for a new organization proposed, the elbow cache. The elbow cache extends the skewed cache organization with a relocation strategy for conflicting blocks.

We show that these skewed designs significantly reduce the conflict problems while consuming up to 56% less dynamic power than a comparably performing 8-way set associative cache. We believe this to be the strongest case in favor of skewed caches presented so far.

Categories and Subject Descriptions: B.3.m [Memory Structures]: Miscellaneous

General Terms: Performance

Keywords: Skewed Caches, Low-Power, Elbow, CAT

1. INTRODUCTION

The emerging trends towards chip multiprocessors and *Symmetric Multi-Threading* (SMT) make shared on-chip caches increasingly common. With multiple threads sharing a cache, it has been shown that the likelihood of destructive sharing increases [10, 21]. To reduce conflict misses, the degree of associativity is often increased. Adding associativity usually comes at a cost in dynamic power consumption. As power consumption has risen to become a first class performance limiting factor, the cost of increasing associativity from an chip power budget perspective may be significant, especially considering highly replicated CMP designs.¹ Consequently there is a need for conflict tolerant caches with low power consumption.

¹Sun's Niagara chip design e.g. contains 16 L1 caches, where each cache is shared by 4 threads [1].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

CF'05, May 4-6, 2005, Ischia, Italy.

Copyright 2005 ACM 1-59593-019-1/05/0005...\$5.00.

Skewed 2-way associative caches have been shown to perform comparable with a 4-way set-associative cache [25]. We will in this paper introduce a new skewed cache design with a replacement strategy based on victim relocation and show that this design outperforms conventional skewed caches. This scheme is called the *elbow cache*.

Since skewed 2-way caches only use half as many bitlines and sense amps as a 4-way associative for each read access, one may suggest that it would consume half as much power. On the other hand it requires two decoders, each connected to half of the wordlines. Because of this the optimized access time and power consumption is different from a conventional 2-way set-associative cache.

We will in this paper present a detailed comparison between the dynamic power consumption of a skewed cache, set-associative caches and the proposed elbow cache organization. The cache designs evaluated was optimized by subarray partitioning making the comparison between organizations as fair as possible.

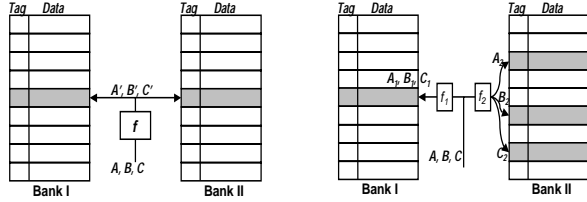
Cache access time is a critical factor in many designs and the additional delay incurred by the skewing function has been viewed as a major drawback of skewed associative caches. We will present a pass-transistor based (XOR) skewing function implementation. By exploiting the early availability of the untranslated bits in the physical address we are able to significantly reduce the delay incurred by the skewing function. We believe this to make skewed caches a much more attractive design choice for timing critical caches.

The contribution of this paper is fourfold:

- We present a power consumption evaluation that shows that skewed caches are viable low power alternatives to highly associative caches.
- We introduce a new cache scheme, the elbow cache, that allows selected victims in a skewed cache to be relocated.
- We describe a fast pass-transistor based skewing function implementation.
- We present various optimizations that could be applied to skewed caches, such as timestamp-based replacement and interleaved way-banks.

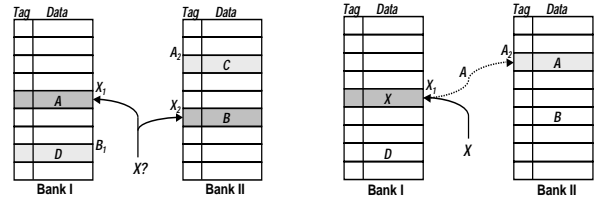
2. THE SKEWED ASSOCIATIVE CACHE

The skewed associative cache was first introduced by Seznec et al. [25, 27]. A skewed cache is conceptually divided into multiple sub-banks each indexed by a different



(a) 2-way set-associative cache

(b) Skewed cache



(a) Elbow cache miss

(b) Elbow relocation

Figure 1: Describing a skewed cache organization.

hash function. The insight behind skewed caches is that cache blocks that map to the same location in one of the banks are likely to map to different locations in the other. This is illustrated in Figure 1 with cache blocks A, B and C conflicting in the set-associative case but not in the skewed case. The effectiveness of the skewed strategy is of course highly dependent on the skewing functions. Thorough discussions on skewing functions are provided in [8, 19, 25, 26, 27]. We have throughout this paper used the XOR-based skewing functions proposed by Bodin and Seznec [8], where two subsets, A_1 and A_2 , of the address bits are XOR'ed together as described below. Let σ be the one bit rotational shift on an n -bit number and \oplus represents a bitwise XOR. The skewing functions are then:

$$f_1(A) = A_1 \oplus A_2 \quad (1)$$

$$f_2(A) = \sigma(A_1) \oplus A_2. \quad (2)$$

3. TIMESTAMP-BASED REPLACEMENT

One of the challenges with skewed caches is the replacement algorithm. Since there are no fixed sets, any combination of victim pairs, one block from each bank, is possible. This makes it difficult to implement an exact ordering-based replacement algorithm like LRU [6]. Instead, approximative algorithms like the *not recently used, enhanced* (NRUE) have been proposed [8]. Another way of providing an ordering between cache blocks is to store a timestamp for each block. At replacement, the block with the least recent timestamp is chosen for eviction. We have found such timestamp-based replacement algorithms to outperform NRUE-replacement for skewed caches at the cost of extra area to hold the timestamps. Results supporting this is presented in Section 6.1.

In [16] Karlsson and Hagersten showed the usefulness of *Cache Allocation Ticks* (CAT)-based timestamps as block survival time metric in the RASCAL-cache. Assuming a uniform reference distribution, the likelihood of a cache line remaining in a cache depends on the number of new allocations performed since the cache line was last touched. Not on the number of cache accesses and/or the number of clock cycles.² This makes the scheme capable of identifying cache lines that are involved in severe conflicts and thus run the risk of being evicted from the cache unexpectedly early.

²The probability for a cache line to remain in a fully associative cache with random replacement is $P_{hit} = (1 - 1/L)^C$, where L is the number of cache lines in the cache and C is the number of cache allocations since the line was last accessed.

Summary of elbow cache actions:

Hit:

- The data and the tag of the two possible locations of X is read.
- The hit is detected
- The current time is written to the timestamp associated with X.

Miss:

- The data and the tag of the two possible locations of X is read.
- X is detected as a cache miss.
- The alternate locations C and D of A and B, are computed and the timestamps of A, B, C and D are retrieved.
- The oldest cache block C is chosen as victim.
- If a secondary victim candidate, C, is selected as victim it resides in a location not accessible to X so another block, A, that is accessible to X, is moved there and X is allocated where A used to be.

Figure 2: Describing an elbow cache replacement.

Another advantage of using cache allocations instead of cache references or cycles is that the resolution of the timestamp is adjusted according to the miss ratio of the cache. This means that the number of bits used in the timestamp can be kept small without the risk of aliasing of the metric. Such CAT-based timestamps give a coarse recency ordering that is still fine-grained enough to distinguish between recent and old data.

The scheme uses a k -bit³ global counter which is incremented each time a new cache block is allocated in the cache. For every cache access, the most significant bits of the current CAT-counter become the new timestamp for the cache block that was accessed. To limit the area overhead of the timestamps only a few of the most significant bits of the CAT-counter value is used. A sensitivity study has shown that 5 bits are sufficient. We have therefore used 5-bit CAT timestamps throughout this paper. Note that the new timestamp simply overwrites the old so no read-modify-write operation is needed.

At the time of a replacement, the timestamps of the cache blocks are retrieved and the CAT-distance d for each block is calculated as shown in Figure 4. T_{curr} is the current value of the global CAT-counter and n is the number of bits in the

³The maximum counter value is four times the number of blocks in the cache ($k = \log_2(blocks) + 2$). This is usually large enough to avoid *aliasing* effects.

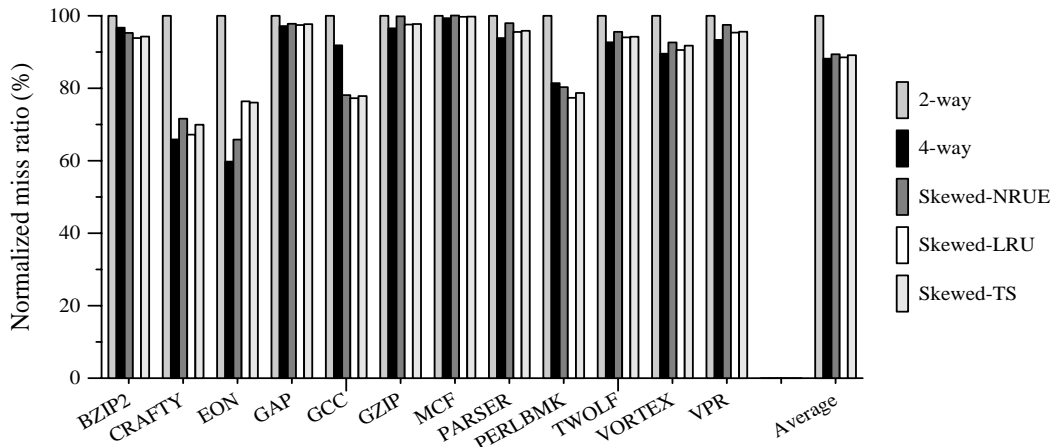


Figure 3: Comparing replacement strategies for skewed caches (SPECint), normalized to a 2-way set-associative cache.

$$d = \begin{cases} T_{curr} - T_{st} & T_{curr} \geq T_{st} \\ T_{curr} + 2^n - T_{st} & T_{curr} < T_{st} \end{cases}$$

Figure 4: Calculating the timestamp distance.

timestamp. The cache block with highest distance is chosen for replacement.

4. THE ELBOW CACHE

The use of timestamps as replacement metric enables a coarse global, temporal ordering of all blocks in the cache. This ordering property is exploited by the scheme proposed in this paper, the elbow cache. An elbow cache extends a skewed organization by carefully selecting its victim and, in the case of a conflict, move the conflicting cache block to its alternate location in the other bank. In a sense, the new data item “uses its elbows” to make space for conflicting data instead of evicting it.

In the common case, a cache hit, the elbow cache works just as a CAT-based skewed cache and updates the timestamp for the accessed block. In the case of a cache miss there are four possible victims: the two *primary* candidates, at the locations hashed to by the address that generated the miss, and also the two *secondary* candidates residing at the alternate locations of the primary candidates. Figure 2a illustrates this, where blocks A and B are the primary candidates (i.e. in the two locations where the new data item X can be placed) and C and D are secondary candidates in the alternate locations of A and B respectively. Among the four blocks, the one with the oldest timestamp is selected as the victim. If a secondary victim candidate is selected, a relocation of the corresponding primary candidate occurs, overwriting the selected victim and freeing space for new data to be filled in (Figure 2b). If the victim is a primary victim candidate no relocation is needed.

In case of a relocation the relocating cache block is written concurrently with the cache fill. This is possible since the fill and relocation always targets different banks. The preceding read of the relocating cache block can also be performed in parallel, in case the selected victim requires a writeback. If the victim data is clean however the relocation read may

compete for cache access with regular accesses. To avoid extra ports for relocation and ensure that relocations don’t stall other cache accesses the cache port arbitration policy may cancel a relocation altogether if no spare cache cycles are found.

4.1 Reducing the number of relocations

A relocation is costly from a power perspective, since a full cache block must be read and written. By imposing restrictions on when to relocate, the power consumption can be reduced. For all the simulation results reported in this paper the following restriction has been applied: First of all, at most one relocation per four misses is allowed.⁴ We found that this restriction prevents extensive power consumption during high miss ratios, with negligible impact on performance. Secondly, we make use of the temporal information provided by our timestamps and only consider blocks with a short reuse distance for relocation. Any other block would be “old” and therefore less likely to be reused in the near future. In our simulations, a relocation can occur only if the selected primary victim has a distance of 3 or less. This significantly reduces the number of relocations without having any substantial impact on hit-ratio performance.

5. METHODOLOGY

The focus of this paper is comparing power consumption for different cache organizations. We have therefore limited our study to cache performance and do not present execution time approximations. Our performance results are reported in terms of miss ratio. Estimations of the dynamic power and access time of the elbow organization was obtained by modifying the Cacti version 3.2 model [28]. We have throughout this paper assumed a process technology of 100nm and simulated caches with a single read/write port. We have also assumed a fixed cycle time for all studied cache organizations, which is a conservative assumption on behalf of the skewed caches since they have shorter access time than the other organizations. This allows us to treat power and energy as exchangeable entities when comparing organizations. The choice of a single read/write port was made

⁴We use a sliding window that allows for up to 16 replacements during the last 64 misses.

since many designs use such as building blocks, through duplication or double-pumping, to build multi-ported caches.

We used the Simics full system simulator [22] to simulate our different workloads. The Simics setup simulated a SPARC-V9 system running an unmodified Solaris 9 operating system. Our evaluation is based on the SPECint 2000 and SPLASH-2 benchmark suites together with two commercial Java workloads. The input data sets for SPLASH was taken from [24] while SPEC was run with the reference dataset.

Our simulated system is a four-way multi-threaded processor with a single data cache shared by all threads. We ran the SPLASH benchmarks with four threads to simulate a single application with several threads sharing the same address space and sharing the same cache. The SPEC benchmarks were used to simulate a machine with several unrelated applications running at once, also sharing the same cache. To create multiprocessor benchmarks from the SPEC suite we mixed four different sub-benchmarks from the suite into a single benchmark. The benchmark mix is shown in Table 3, Appendix A. The two simulated Java middleware workloads, ECperf and SPECjbb, was described in detail by Karlsson et al. [17]. The Java workloads results are from uniprocessor simulations.

The SPLASH benchmarks were studied only in the parallel region of the code, and the caches were all warmed up before the measurements started. For the SPEC runs we fast-forwarded 10 Billion instructions into the benchmark before warming up the caches for 500 Million instructions. The measurements were then obtained for the following 2 Billion instructions.

The skewed and elbow caches are compared to a number of conventional set-associative caches with LRU replacement and various degrees of associativity. We present results for a 32KB data cache with 64-byte blocks. However, we have observed similar results for 16KB and 64KB caches.

6. CACHE PERFORMANCE RESULTS

In this section we present miss ratio performance results for a 2-way skewed cache, a 2-way elbow cache as well as 4- and 8-way set-associative caches. We present our performance evaluation in the terms of miss ratio normalized to a 2-way set-associative cache. Table 2, Appendix A shows the absolute miss ratios for the 2-way set-associative cache.

6.1 CAT-Timestamp Metric Performance

Our first performance experiment is a comparison between different replacement metrics for skewed caches. We model a simple uniprocessor system and run the SPECint benchmarks. We simulate a perfect LRU replacement, a NRUE replacement and replacement based on CAT-timestamps (TS). The results are shown, as miss ratios normalized to a two-way set-associative cache, in Figure 3 where we also included a four-way set-associative cache for comparison. The results show that all replacement strategies perform similarly, but that timestamp replacement performs slightly better than NRUE-replacement.

6.2 Shared Cache Performance

The results of the SPLASH simulations are shown in Figure 5 When increasing associativity from 2-way to 4-way, we observe a substantial miss ratio reduction. The reduction beyond 4-way is very limited for most benchmarks. We

do note however, that both the skewed and the elbow cache perform well in comparison to the other two caches. For the commercial java workloads we find that, in the case of ECperf, all organizations except the 2-way perform almost identically. In SpecJBB however the elbow cache performs better than both the 4-way and the skewed cache.

The SPEC benchmarks results are shown in Figure 6. Unlike the SPLASH benchmarks, the reductions from 2-way set-associativity are quite small, around 10–20%. On the other hand, the 2-way set-associative cache has significantly higher miss ratio to begin with (See Table 2, Appendix A). Like in the SPLASH case above, 4-way set-associativity captures most of the conflicting data. Benchmarks 1, 3 and 7 (BMK_01, BMK_03 and BMK_07) are exceptions however, and show a much higher reduction for the 8-way, skewed and elbow caches.⁵ This is indicating a high amount of conflict misses and that conflict reduction beyond 4-way set-associativity is worthwhile. The results also show that both skewed and our proposed elbow cache are effective for these benchmarks.

7. POWER ESTIMATES

The power consumption of current architectures consists almost entirely of dynamic power. However as process technology advances the static power consumption is projected to make up an increasingly large portion of the total power. In this section we compare dynamic power consumption based on energy estimations for the skewed cache, the elbow cache and the 2, 4 and 8-way set-associative caches. We assume that the static power is similar for all of these designs.

7.1 Dynamic Power

As process technology scale to smaller feature sizes, the power consumed by the cache’s sense-amplifiers makes up an increasingly large fraction of the total power consumed by a cache read. This can be observed by running Cacti with different technology parameters. Since the number of words read out in parallel determines how many bit-lines and sense-amplifiers that need to be activated during an access, associativity comes at an increasing cost in power.

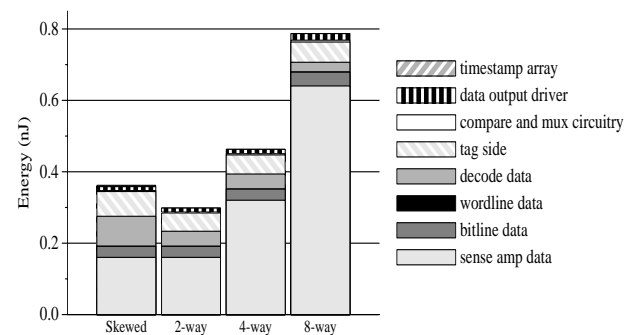


Figure 7: Energy breakdown per hit for a 32KB cache.

⁵The higher miss ratios of the 4-way cache for BMK_01 and BMK_03 might seem unintuitive but results from non-optimal replacement (LRU).

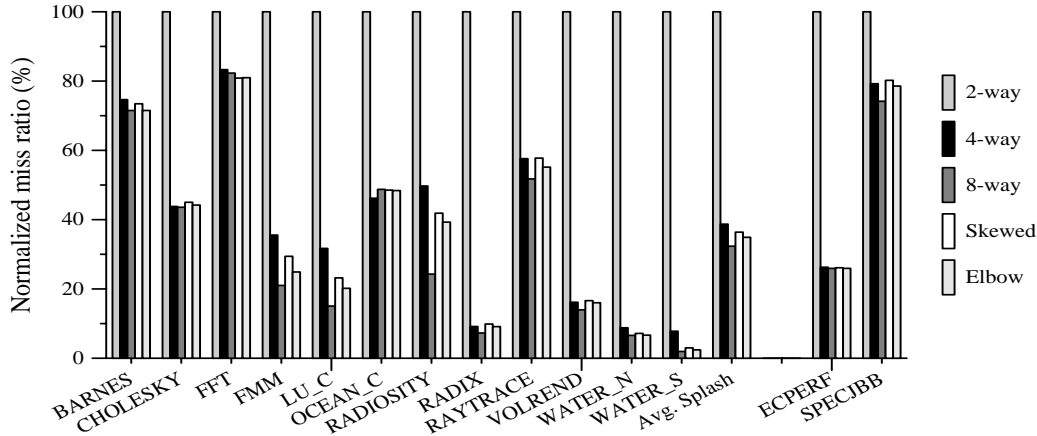


Figure 5: Normalized miss ratio for one multi-threaded application sharing one data cache (SPLASH), compared with a 2-way set-associative cache (miss ratios listed in Appendix A).

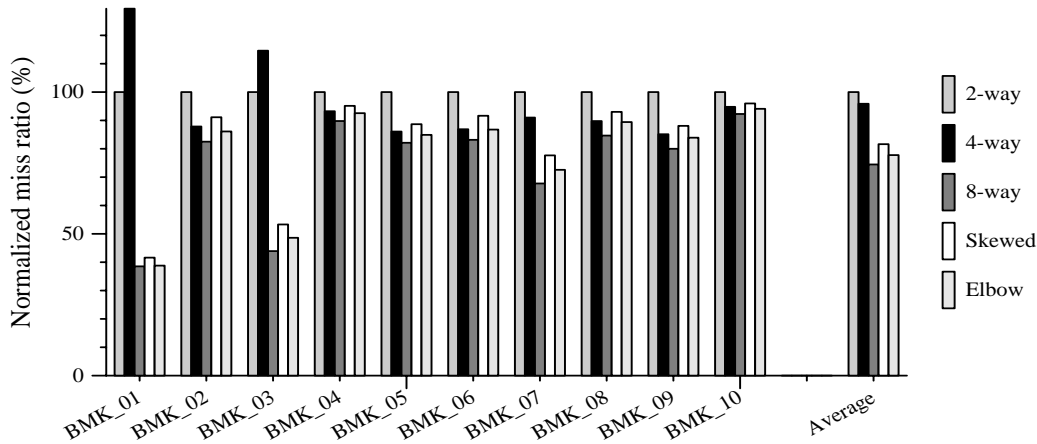


Figure 6: Normalized miss ratio for four single-threaded applications sharing one data cache (SPECint), compared with a 2-way set-associative cache (miss ratios listed in Appendix A).

A common approach to reduce the access time and power consumption of a cache is to divide the SRAM-cell array into sub-arrays. Depending on the optimal aspect-ratio and sub-array size, bit-line-wise and word-line-wise divisions are applied. Hence an 8-way associative cache is not necessarily twice as wide as a 4-way associative cache, nor consumes twice as much power. In our evaluation, the sub-array divisions used for each particular organization was given by Cacti’s optimization function. This function takes energy as well as access time, area and aspect ratio into account.

As can be seen in Figure 7, the per lookup energy for a cache hit, in the skewed/elbow cache⁶ is 17% higher than for a 2-way set-associative cache, but only 75% and 44% of the energy consumption of a 4-way and 8-way cache respectively. The majority of this difference comes from the data sense amps. For the skewed cache the energy cost of having two decoders can also be seen in Figure 7.

The dynamic power consumed by a cache can be computed by equation (3) below, where P_{Ld} denotes the power

⁶The power estimates are identical assuming that both caches use timestamp-based replacement.

consumed by loads and P_{St} represents the power associated with stores:

$$P_{Dynamic} = P_{Ld} + P_{St} \quad (3)$$

$$P_{LdConv} = P_{Hit} + Ld_{MissFrac} * P_{Fill} \quad (4)$$

$$P_{LdElbow} = P_{LdConv} + Ld_{MissFrac} * P_{Relocation} * freq_{Reloc} \quad (5)$$

To simplify our study we have opted to exclude store power consumption.⁷ The P_{Fill} term is the additional power consumed while reading a new block from the L2 cache and filling it into the L1 cache when a cache miss occurs. The total dynamic load power for a conventional (or skewed) cache can be expressed as (4). More details of our implementation assumptions are presented in Section 9.

7.2 Elbow Relocation Power

Due to the energy associated with a relocation in the elbow cache, the cache power consumption is also dependent on how often relocations occur. This, in turn, depends on

⁷Since the tag and data is looked up sequentially in the case of a write, the write power is approximately the same for all organizations.

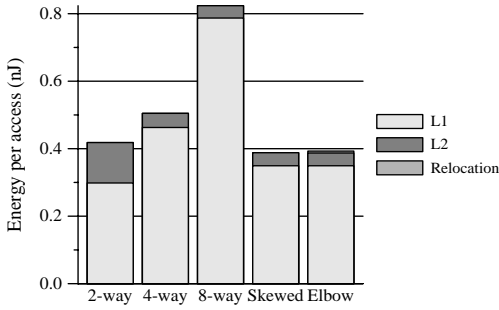


Figure 8: Average energy consumption per load for different cache architectures across all SPLASH benchmarks.

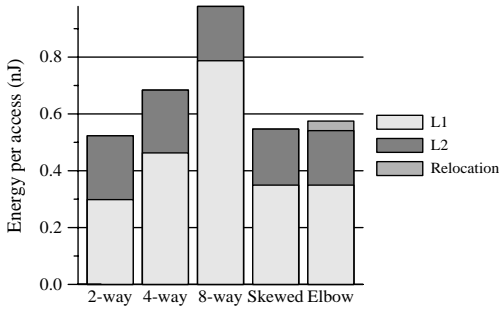


Figure 9: Average energy consumption per load for different cache architectures across all SPEC benchmark mixes.

the miss ratio. Hence, in the case of zero percent miss ratio and therefore no relocations, the elbow cache will consume the same amount of power as a skewed cache. For non-zero miss rates the power consumption will increase linearly with the miss rate and relocation frequency. Formula (5) shows the average load power for the elbow cache. $freq_{Reloc}$ is the probability of a relocation whenever a miss occurs.

7.3 Memory System Power

When evaluating power consumption for a certain level of the memory hierarchy it is also important to consider the power cost of accesses to the next level. To account for this we have included the power consumption of a 64-byte block load from an 8-way 1MB level two (L2) cache, with serial tag and data array lookups, in the fill power term (P_{Fill}). The traffic between a large L2 cache and the lower levels of the memory hierarchy should be similar for all the different L1 cache systems. The power consumption of other parts of the memory hierarchy, like the bus and memory, will therefore not affect our evaluation and is excluded.

Figures 8 and 9 show the average energy-per-access of the SPLASH and SPEC simulations broken down into level one and level two cache accesses. The energy consumption for relocations in the elbow cache is also included. As can be seen in the graphs, the total power consumption is highly dependent on the miss ratio. For example, the total cache system power consumption for SPLASH is higher for the 2-way associative cache than for the skewed and elbow caches, although it is lower from an L1 power perspective.

8. SKEWING DELAY AND ACCESS TIME

The access time of a cache is very important for system performance, since it may affect the processor cycle time [12, 13]. Timing wise there are a few issues surrounding skewed caches that one must bear in mind. Since the elbow cache is an extension of the skewed cache concept it too will inherit these issues. One of the drawbacks of skewed caches, primarily affecting L1 caches, is that an additional number of physical address bits are required by the skewing functions to calculate the index. This makes virtually-indexed caches infeasible for realistic cache sizes. Another perhaps more serious drawback is the additional critical-path delay introduced by skewing functions. This extra delay can be substantial in the context of timing-critical L1 caches.

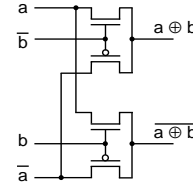


Figure 10: Pass-transistor based skewing function.

To reduce the skewing delay we use a pass-transistor layout [31] for the XOR-functions (Figure 10). If bits b and \bar{b} in the address are available before a and \bar{a} , this design can be used to reduce the skew-delay to a fraction of that of a logical gate. If the bit b is taken from the part of the address that does not need to be translated via the TLB, this criteria is fulfilled. This solution has limited the number of bits that are available to the skewing function. If we assume a page size of 8KB,⁸ an address A can be written as $A = \{a_N, \dots, a_0, b_{12}, \dots, b_0\}$ where bits b_n belong to the physical part and a_n to the virtual part. Bits $\{b_5, \dots, b_0\}$ are used to index the word in each block. For our 32KB 2-way elbow cache that has 256 blocks in each bank, bits $\{a_0, b_{12}, \dots, b_6\}$ would be XOR'ed with bits $\{a_9, \dots, a_1\}$ to produce an ideal skewing function.⁹ If we instead only apply the XOR-function to the sets $\{b_{12}, \dots, b_6\}$ and $\{a_8, \dots, a_1\}$, and let a_0 fill in the missing most significant bit, we get a more restricted skew, but one were the pass-transistor layout is possible. All performance results reported in this paper use this type of skewing function, since the performance of skewed caches is not very sensitive to this type of restriction [30]. In our case, the pass-transistor design yield an access time overhead for the skewing function of less than 0.7%. This is however dependant on the width of the pass-transistors.

By changing the sub-array division of the cache, speed can be traded for power and vice versa. We have chosen a configuration that enables our skewed/elbow cache to have a shorter access time, including skewing delay, than that of the 4-way or 8-way set-associative caches. Table 1, lists the access time and energy consumption results, together with the sub-array configurations used. Our Cacti results yield a skewed cache with an access time that is approximately the same as the 4-way set-associative cache and 12% lower than

⁸Including page coloring.

⁹The inverse address bits, also needed by the address decoder, are obtained in the same way.

Arch.	Access Time [ns]	Energy [nJ]	Sub-array configuration Ndwl-Ndbl-Nspd-Ntwl-Ntbl-Ntspd
2-way	0.744	0.298	4-2-1-1-2-2
4-way	0.741	0.463	8-1-1-1-2-1
8-way	0.830	0.787	8-1-1-1-2-1
Skewed	0.740	0.349	4-2-1-1-2-2

Table 1: Access time, load energy consumption and subarray partitioning for the evaluated caches.

the 8-way. Note that this study focuses on timing-critical caches where tag and data are accessed in parallel. If the latency of sequential tag and data lookup is tolerable, such a design is of course far better from a dynamic power perspective. This may however result in up to 60% higher access time [32]. We observe that the 2-way associative cache has a higher access time than the Skewed cache despite the additional skewing delay. This is a result of the fact that in the skewed cache design only half as many bit cells are connected to each decoder, and the shorter word line delay more than compensates for the skewing delay.

9. IMPLEMENTATION DETAILS

This section describes the assumed implementation of the skewed/elbow cache on which our modified Cacti-model was based.

First we have added a separate timestamp structure to hold the timestamp of each cache block in the cache. Assuming a 64-byte cache block size and 5-bit timestamps this gives an extra area cost of approximately 2% of the data array. The timestamps are separated per bank so that the first bank has one timestamp array and the second has another.

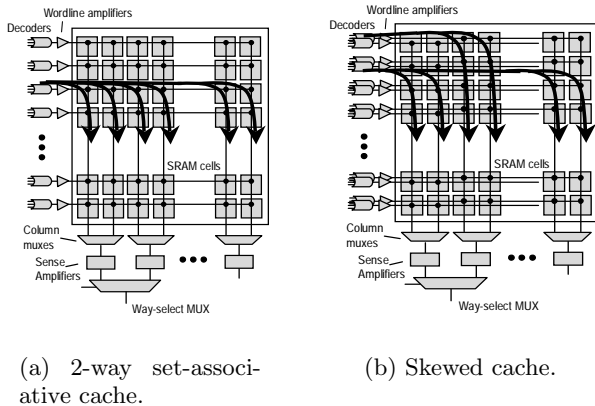


Figure 11: Way interleaving the logical banks of the skewed cache.

Secondly we have used a physical layout of our skewed cache such that the two logically separate way-banks are interleaved into one physical bank. Therefore bit-line n from the first bank is placed close to bit-line n from the second bank, similar to the layout that of a 2-way set-associative array (Figure 11). By doing so, we get a better aspect-ratio on the array and also simplify the design of the way-select logic. This helps when implementing the relocation

functionality in the elbow cache, since the bit-lines of any two bit-cell pairs affected by the move will be physically adjacent.

To be able to look up both logical banks concurrently, two decoders are required. Each decoder is connected to half of the cells in the array. We assume that the additional word-lines can be routed in a different metal layer making extra word-line spacing unnecessary. It is still possible to make sub-array divisions with this arrangement, but when splitting an array bit-line-wise, restrictions must be made on the skewing functions so that an address is always mapped into the same sub-array by both way-functions. In the particular layout for our 32KB elbow cache however, these restrictions are less rigid than the ones already imposed by our solution to reduce the skewing delay, and do therefore not affect performance.

10. RELATED WORK

Some of the earlier work that has addressed the issue of reducing cache conflicts [13] is presented here. Jouppi [29] presented the *victim cache*, primarily for use in direct-mapped caches. Topham and Gonzales [15] studied the use of hash-functions for indexing a cache. Agarwal and Pudar [3] suggested *column-associativity* for improving direct-mapped caches. Seznec and Bodin [27] pioneered the work on *skewed-associative* caches. Further work on skewed caches is presented in some of their later papers [8, 25, 26].

The power consumption of current architectures consist almost entirely of dynamic power. However as process technology advances the static power consumption is projected to make up an increasingly large portion of the total power. The Elbow cache organization presented in this study is evaluated from a dynamic power consumption aspect. The static power is expected to be the same as for a conventional 2-way associative design. However the proposed architecture may very well be combined with leakage power reduction techniques [11, 18, 33].

Recently many proposals have been made of how to reduce the dynamic power consumption in various on-chip memory structures [2, 5, 7, 9, 14, 20, 34, 35].

To overcome the latency penalty of a sequential access, Powell et al. proposes *way-prediction* [23], to predict which way to access. This method has the drawback of multiple hit times, a correctly predicted hit and a mispredicted hit, and a significant prediction table size. An interesting alternative to consider is to apply skewed caches in this context. A skewed cache of lower associativity, and hence fewer ways to choose from, could be used to simplify such a design and improve the correct way-prediction rate.

Another proposed power saving approach is to shutdown unnecessary parts of the cache during periods of low utilization. Albonesi proposed the idea of Selective cache ways [4] where unneeded ways in a cache are shutdown.

11. CONCLUSION

Skewed caches are often very effective in reducing conflicts. The multi-threaded nature of emerging processors and the shared caches that follow, are likely to increase the need for such highly conflict tolerant designs. In this paper we revisit the skewed caches and evaluate them in the context of power consumption. A timestamp-based replacement metric is proposed.

We also introduce a new cache organization called the elbow cache. It extends a skewed cache by relocating conflicting cache blocks to their alternate locations. In this study we have evaluated the elbow cache in the context of data caches. The idea of relocation of data to reduce conflicts can however be applicable in many other areas as well, such as instruction caches, TLBs, BTBs etc.

Dynamic power evaluations for both the skewed and elbow caches are presented and compared to traditional set-associative caches. We find that despite the extra decoder and timestamp structure, the skewed caches per access power consumption is lower than that of an 8-way set-associative cache. When also taking L2 power into account, we find that the skewed and elbow cache configurations consume significantly less total power than both 4-way and 8-way set-associative cache configurations.

From our miss ratio evaluation we conclude that the timestamp based skewed cache shows surprisingly good performance compared with set-associative caches. The elbow cache further reduces the miss ratio at the cost of more complexity. Considering the constrained power budgets of today's processors, a skewed or an elbow cache should make an interesting alternative to those more traditional designs.

12. ACKNOWLEDGMENTS

We would like to thank Mark Hill for contributing to the initial discussion leading up to the idea of relocation. We would also like to thank Magnus Ekman for help with Cacti. This work is funded by the PAMP research program, supported by the Swedish Foundation for Strategic Research.

APPENDIX

A. TABLES

SPECint mix		SPLASH-2 & Java	
BMK_01	4.91%	BARNES	2.49%
BMK_02	6.74%	CHOLESKY	5.75%
BMK_03	5.70%	FFT	6.62%
BMK_04	18.99%	FMM	2.21%
BMK_05	10.09%	LU_C	2.96%
BMK_06	8.95%	OCEAN_C	5.58%
BMK_07	9.26%	RADIOSTY	3.94%
BMK_08	10.44%	RADIX	15.40%
BMK_09	9.60%	RAYTRACE	5.06%
BMK_10	20.79%	VOLREND	6.95%
		WATER_N	4.97%
		WATER_S	5.57%
Avg. SPEC	10.55%	Avg. Splash	5.62%
		ECperf	6.09%
		SpecJBB	10.26%

Table 2: Miss ratios for the 32KB, 2-way set-associative reference cache.

B. REFERENCES

[1] Sun's Niagara Pours on the Cores. Microprocessor Report Newsletter, September 2004.

	SPECint mix		SPECint mix
BMK_01	BZIP2_GFX GAP GZIP_LOG PERLBMK_M	BMK_06	GCC_166 GZIP_LOG PERLBMK_P VORTEX_BE2
BMK_02	VORTEX_BE1 BZIP2_SRC GCC_200 GZIP_RND	BMK_07	CRAFTY GCC_EXPR0 GZIP_SRC PERLBMK_SM2
BMK_03	PERLBMK_SM1 VORTEX_BE3 EON_COOK GCC_INTEGR	BMK_08	VPR_PLAC EON_KAJI GCC_SCILAB PARSER
BMK_04	MCF PERLBMK_SM3 VPR_ROUT EON_RUSH	BMK_09	PERLBMK_SM4 BZIP2_PGR GCC_200 GZIP_SRC
BMK_05	GZIP_GFX PERLBMK_D TWOLF BZIP2_PGR	BMK_10	GAP MCF TWOLF CRAFTY

Table 3: The mix of four SPEC sub-benchmarks used to create the four-cpu multiprocessor benchmarks.

[2] J. Abella and A. Gonzalez. Power Efficient Data Cache Designs. In *Proceedings of the 21st International Conference in Computer Design*, 2003.

[3] A. Agarwal and S. D. Pudar. Column-Associative Caches: A Technique for Reducing the Miss Rate of Direct-Mapped Caches. In *Proceedings of the 20th International Symposium on Computer Architecture*, pages 179–190, May 1993.

[4] David H. Albonesi. Selective Cache Ways: On-Demand Cache Resource Allocation. In *International Symposium on Microarchitecture*, pages 248–, 1999.

[5] B. Bannon and T. N. Vijaykumar. Reactive-Associative Caches. In *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques*, 2001.

[6] L. A. Belady. A study of replacement algorithms for a virtual storage computer. *IBM Systems Journal*, 5:78–101, 1966.

[7] N. Bellas, I. Hajj, and C. Polychronopoulos. Using dynamic management techniques to reduce energy in high-performance processors. In *Proceedings of the 1999 International Symposium on Low Power Electronics and Design (ISLPED)*, 1999.

[8] F. Bodin and A. Sez nec. Skewed associativity improves program performance and enhances predictability. In *IEEE Transactions on Computers*, May 1997.

[9] Magnus Ekman, Fredrik Dahlgren, and Per Stenström. TLB and Snoop Energy-Reduction using Virtual Caches in Low-Power Chip-Multiprocessors. In *Proceedings of the 2002 International Symposium on Low Power Electronics and Design (ISLPED)*, 2002.

[10] Andrew Erlichson, Basem A. Nayfeh, Jaswinder Pal Singh, and Kunle Olukotun. The Benefits of Clustering in Shared Address Space Multiprocessors: An Applications-Driven Investigation. In *Supercomputing*, 1995.

[11] K. Flautner, N. Kim, S. Martin, D. Blaauw, and

- T. Mudge. Drowsy Caches: Simple Techniques for Reducing Leakage Power. In *Proceedings of the 29th Annual International Symposium on Computer Architecture (ISCA'02)*, 2002.
- [12] R. Heald, K. Shin, V. Reddy, I.-F. Kao, M. Khan, W. L. Lynch, G. Lauterbach, and J. Petolino. 64kB Sum-Addressed-Memory Cache with 1.6ns Cycle and 2.6ns Latency. *IEEE Journal of Solid-State Circuits* 33, page 16821689, 1998.
- [13] M.D. Hill. *Aspects of Cache Memory and Instruction Buffer Performance*. PhD thesis, University of California, Berkeley, 1987.
- [14] J. Jalminger and P. Stenström. Improvements of Energy-Efficiency in Off-Chip Caches by Selective Prefetching. *Microprocessors and Microsystems*, 2001.
- [15] Nigel P. Topham and Antonio Gonzalez. Randomized Cache Placement for Eliminating Conflicts. *IEEE Transactions on Computers*, 48(2):185–192, 1999.
- [16] M. Karlsson and E. Hagersten. Timestamp-based Selective Cache Allocation. In *Proceedings of the Workshop on Memory Performance Issues*, June 2001. held in conjunction with the 28th International Symposium on Computer Architecture (ISCA28).
- [17] M. Karlsson, K. Moore, E. Hagersten, and D. A. Wood. Memory System Behavior of Java-Based Middleware. In *Proceedings of the Ninth International Symposium on High Performance Computer Architecture (HPCA-9)*, Anaheim, California, USA, February 2003.
- [18] Stefanos Kaxiras, Zhigang Hu, and Margaret Martonosi. Cache decay: Exploiting generational behavior to reduce cache leakage power. In *Proceedings of the 28th International Symposium on Computer Architecture*, pages 240–251, 2001.
- [19] Mazen Kharbutli, Keith Irwin, Yan Solihin, and Jaejin Lee. Using Prime Numbers for Cache Indexing to Eliminate Conflict Misses. In *Proceedings of the 10th International Symposium on High-Performance Computer Architecture*, pages 288–299, 2004.
- [20] J. Kin, M. Gupta, and W. H. Mangione-Smith. The filter cache: An energy efficient memory structure. In *Proceedings of the 30th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO-30)*, 1997.
- [21] Jack L. Lo, Luiz Andre Barroso, Susan J. Eggers, Kourosh Gharachorloo, Henry M. Levy, and Sujay S. Parekh. An Analysis of Database Workload Performance on Simultaneous Multithreaded Processors. In *Proceedings of the 25th Annual International Symposium on Computer Architecture (ISCA'98)*, pages 39–50, 1998.
- [22] P. S. Magnusson, M. Christensson, D. Forsgren, J. Eskilson, G. Hällberg, J. Högberg, A. Moestedt, F. Larsson, and B. Werner. Simics: A Full System Simulation Platform. *IEEE Computer*, February 2002.
- [23] M. D. Powell, A. Agarwal, T. N. Vijaykumar, Babak Falsafi, and Kaushik Roy. Reducing Set-Associative Cache Energy via Way-Prediction and Selective Direct-Mapping. In *Proceedings of the 34th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO-34)*, 2001.
- [24] S. Woo, M. Ohara, E. Toorie, J.P. Singh, and A. Gupta. The SPLASH-2 Programs: Characterization and Methodological Considerations. In *Proceedings of the 22th International Symposium on Computer Architecture*, pages 24–36, June 1995.
- [25] A. Seznec. A case for two-way skewed associative caches. In *Proceedings of the 20th International Symposium on Computer Architecture*, pages 169–178, May 1993.
- [26] A. Seznec. A new case for skewed-associativity. Internal Publication No 1114, IRISA-INRIA, July 1997.
- [27] A. Seznec and F. Bodin. Skewed-associative caches. In *Proceedings of PARLE '93, Munich*, pages 305–316, June 1993.
- [28] P. Shivakumar and N. Jouppi. CACTI 3.0 An integrated Cache Timing, Power and Area Model. Technical Report 2001/2, DEC Western Research Lab, 2001.
- [29] N. P. Jouppi. Improving Direct-Mapped Cache Performance by the addition of a Small Fully-Associative Cache and Prefetch Buffers. In *Proceedings of the 17th International Symposium on Computer Architecture*, June 1990.
- [30] Hans Vandierendonck and Koen De Bosschere. Trade-offs for Skewed-Associative Caches. In *Parallel Computing (PARCO)*, September 2003.
- [31] N. H. E. Weste and K. Eshraghian. *Principles of CMOS VLSI Design*. Addison-Wesley, second edition, 1993.
- [32] S. Wilton and N. Jouppi. An enhanced access and cycle time model for on-chip caches, 1994.
- [33] Se-Hyun Yang, Michael D. Powell, Babak Falsafi, Kaushik Roy, and T. N. Vijaykumar. An Integrated Circuit/Architecture Approach to Reducing Leakage in Deep-Submicron High-Performance I-Caches. In *International Symposium on High-Performance Computer Architecture (HPCA)*, 2001.
- [34] C. Zhang, F. Vahid, and W. Najjar. A highly configurable cache architecture for embedded systems. In *Proceedings of the 30th Annual International Symposium on Computer Architecture (ISCA'03)*, 2003.
- [35] Michael Zhang and Krste Asanovic. Highly-Associative Caches for Low-Power Processors. In *Proceedings of Kool Chips Workshop held in conjunction with International Symposium on Microarchitecture (MICRO-33)*, Monterey, CA, December, 2000.