

Eliminating Cache Conflict Misses Through XOR-Based Placement Functions

Antonio González*, Mateo Valero*, Nigel Topham[†] and Joan M. Parcerisa*

* Departament d'Arquitectura de Computadors
Universitat Politècnica de Catalunya
c/ Jordi Girona 1-3, 08034 Barcelona (Spain)

Email: {antonio,mateo,jmanel}@ac.upc.es

[†] Department of Computer Science
University of Edinburgh
JCMB, Kings Buildings, Edinburgh (UK)

Email: npt@dcs.ed.ac.uk

Abstract

This paper makes the case for the use of XOR-based placement functions for cache memories. It shows that these XOR-mapping schemes can eliminate many conflict misses for direct-mapped and victim caches and practically all of them for (pseudo) two-way associative organizations. The paper evaluates the performance of XOR-mapping schemes for a number of different cache organizations: direct-mapped, set-associative, victim, hash-rehash, column-associative and skewed-associative. It also proposes novel replacement policies for some of these cache organizations. In particular, it presents a low-cost implementation of a LRU replacement policy which demonstrates a significant improvement over the pseudo-LRU replacement previously proposed. The paper shows that for a 8 Kbyte data cache, XOR-mapping schemes approximately halve the miss ratio for two-way associative and column-associative organizations. Skewed-associative caches, which already make use of XOR-mapping functions, can benefit from the LRU replacement and also from the use of more sophisticated mapping functions. For two-way associative, column-associative and two-way skewed-associative organizations, XOR-mapping schemes achieve a miss ratio that is not higher than 1.10 times that of a fully-associative cache. XOR mapping schemes also provide a very significant reduction in the miss ratio for the other cache organizations, including the direct-mapped cache. Ultimately, the conclusion of this study is that XOR-based placement functions unequivocally provide highly significant performance benefits to most cache organizations.

Keywords: cache memory, XOR-based placement functions, conflict misses.

1 Introduction

The use of XOR functions to map memory addresses onto a set of memory modules has been studied extensively in the last decade; for example, see [8], [15], [21], [11], [16], [10], [17] and [23]. It has proven to be an effective way to distribute memory addresses to memory modules in a pseudo-random way. In that context, the aim is to allow multiple memory references to proceed in parallel by maximizing the probability that they will access different memory modules. The effect of random distribution can be also beneficial for cache memories if it is used to map memory addresses onto

cache data sets. In this case, the desired effect would be the removal of conflict misses. In fact, a cache memory with a pure random placement would be equivalent in terms of hit ratio to a fully-associative cache with a random replacement. This paper shows that this can be achieved with simple XOR-mapping schemes together with a (pseudo) two-way associative organization. For direct-mapped and victim caches, the reduction in number of conflict misses is also very high although they are not completely removed.

Despite the potential benefits of XOR-mapping schemes, there are very few proposals in the literature using these schemes for cache memories. The most notable are the skewed-associative cache [18] [19] and the cache memory of the HP 7100 [5].

In this paper, we present a study of the use of XOR-mapping schemes on a number of different cache organizations: direct-mapped [20], set-associative [20], victim [14], hash-rehash [3], column-associative [4] and skewed-associative [19]. The paper also proposes a low-cost implementation of the LRU replacement policy for use with XOR-mapping functions. It is shown that this replacement policy provides a significant improvement for the column-associative and the two-way skewed-associative cache.

Two different types of XOR-mapping schemes have been evaluated: a simple bitwise XOR of two fields of the address and the polynomial mapping proposed by B. Rau [17].

For the bitwise XOR scheme, a column-associative cache with LRU replacement, and without swapping, has the lowest miss ratio. This miss ratio is significantly lower than that of a four-way set associative cache and very close to that of a fully-associative cache. A two-way associative cache with an XOR-mapping function yields almost the same hit ratio. We found that a two-way skewed-associative cache has a significantly higher miss ratio when it uses the replacement policy originally proposed by its author. This miss ratio is about the same as that of a victim cache with a bitwise XOR-mapping. However, when using the LRU replacement policy proposed in this paper, the skewed-associative cache achieves a miss ratio very similar to the column-associative and the two-way associative organizations. When swapping is incorporated in a column-associative cache, the overall miss ratio increases slightly, due to the use of the XOR-mapping. However, in this case most of the hits are obtained with a single probe, which may reduce the average access time. A direct-mapped cache exhibits the highest miss ratio, but even in this case, an XOR-mapping function yields very significant improvement.

Polynomial mapping provides marginal advantages for the column-associative and for the two-way associative organizations. However, it is more effective for the two-way skewed associative cache. With this type of mapping, the two-way skewed-associative cache achieves the lowest miss ratio, which is practically identical to that of a fully associative cache (0.8% higher).

Overall, the two-way skewed-associative, column-associative and two-way associative organizations exhibit a similar miss ratio. Miss ratio is not the only parameter to consider when evaluating a cache memory. The most relevant performance metric is the average memory access time, which depends on the access time of the cache memory, the miss ratio and the miss penalty. Genuinely set-associative caches have higher hit times than pseudo-associative caches, though the latter may require two probes to detect a hit. In addition, LRU replacement requires more hardware for the column-associative and skewed-associative organizations. Thus, the most effective organization in practice will depend on the hardware implementation.

The rest of this paper is organized as follows. Section 2 summarizes related work. Some basic concepts are reviewed in section 3, which also describes the evaluation methodology. The performance of conventional mapping functions is evaluated in section 4 for a selection of cache configurations. Section 5 explores the benefits of using XOR-mapping functions in those cache organizations. The implementation of LRU replacement in the presence of an XOR mapping is discussed and evaluated in section 6. The effectiveness of polynomial mapping is analyzed in section 7. Section 8 evaluates the effect of swapping in the column-associative cache. Finally, the main conclusions of this work are summarized in section 9.

2 Related work

There are remarkably few papers on the use of alternative mapping schemes for cache memories. The first computers based on the HP Precision Architecture Processor [7] made use of XOR-mapping functions in order to index the TLB. In these machines, the 11-bit TLB index was obtained by the exclusive OR of two 9-bit fields, one from the virtual page number and the other from the space ID, appended to two other bits of the space ID. Earlier machines that used a XOR-mapping function to index the TLB were the IBM 3033 [13] and the Amdahl 470 [1].

The use of XOR-mapping schemes in order to obtain a pseudo-random placement has been suggested by other authors as reported in [20]. In [20], a comparison of a pseudo-random placement against a set-associative one was performed. It concluded that random mapping had a small advantage in most cases, but that the advantage was not significant. We will show in this paper that for current workloads and cache organizations, this advantage can be very large.

Hashing the process ID with the address bits in order to index the cache memory was evaluated in [2] for a multiprogrammed environment. Results were provided for just one trace, which shown that this scheme could reduce the miss ratio.

In practical systems, like the HP PA 7100, limited and undocumented use of XOR-mapping schemes has occurred, but there is currently no established body of published results analyzing the true benefits of alternative mapping schemes.

More recently, the use of XOR-mapping functions was proposed in skewed-associative caches [18] [19]. A two-way skewed-associative cache consists of two banks of the same size that are accessed simultaneously with two different hashing functions. In that paper, a family of mapping functions was defined as follows. Assume that the cache memory consists of a 2^l lines of 2^b bytes each. A memory address $A = \langle a_{n-1}, a_{n-2}, \dots, a_0 \rangle$ comprises the following fields: $A = \langle A_3, A_2, A_1, A_0 \rangle$ such that $A_0 = \langle a_{b-1}, \dots, a_0 \rangle$; $A_1 = \langle a_{1+b-2}, \dots, a_b \rangle$; $A_2 = \langle a_{2l+b-3}, \dots, a_{l+b-1} \rangle$; and $A_3 = \langle a_{n-1}, \dots, a_{2l+b-2} \rangle$. Let \oplus denote the bitwise exclusive OR; let \bullet denote the bitwise AND operation and let T be any $(l-1)$ -bit number (a good choice for T would be 1010...10); let $\bar{T} = 2^{l-1} - T$. The family of twin XOR-based placement functions are defined as:

$$f_0^T: \{0 \dots 2^n - 1\} \rightarrow \{0 \dots 2^{l-1} - 1\}$$

$$A = \langle A_3, A_2, A_1, A_0 \rangle \rightarrow ((A_2 \bullet T) \oplus A_1, A_0)$$

$$f_1^T: \{0 \dots 2^n - 1\} \rightarrow \{0 \dots 2^{l-1} - 1\}$$

$$A = \langle A_3, A_2, A_1, A_0 \rangle \rightarrow ((A_2 \bullet \bar{T}) \oplus A_1, A_0)$$

In [18], it was proposed a pseudo-LRU replacement by associating a one-bit flag to each line in bank 0. If the requested data is found in bank 0, the corresponding line flag is set, whereas it is reset if the data is found in bank 1. On a miss, the flag of the line selected in bank 0 is read and its value determines the bank where the missing data is to be placed.

Using a different workload from that used in this paper, it was observed that the miss ratio of the two-way skewed-associative cache was lower than that of a victim cache (with four lines in the victim buffer) and similar to the miss ratio of a four-way set associative cache.

3 Preliminaries

Whenever a line of main memory is brought into cache a decision must be made on which line, or set of lines, in the cache will be candidates for storing that memory line. This **line placement** policy is one of the least researched aspects of cache design. Direct-mapped caches typically extract a field of l bits from the address and use this to select one line from a set of 2^l . Whilst simple, and trivial to implement, this mapping function is not robust. The principal weakness of this function is its susceptibility to repetitive conflict misses. For example, if C is the cache capacity and B is the line size, then addresses a_1 and a_2 map to the same cache line if $\lfloor a_1/B \rfloor \bmod C = \lfloor a_2/B \rfloor \bmod C$. If a_1 and a_2 map to the same cache line, then addresses $a_1 + k$ and $a_2 + k$ are guaranteed to also map to identical cache lines, for any integer $k \geq B$. There are two common cases when this happens:

- when accessing a stream of addresses $A = \{a_0, a_1, \dots, a_m\}$ if a_i collides with a_{i+k} , then there may be up to $(m-k)$ conflict misses in this stream.
- when accessing elements of two distinct arrays b_0 and b_1 , if $b_0[i]$ collides with $b_1[j]$ then $b_0[i+k]$ will collide with $b_1[j+k]$, for any integer $k \geq B$.

w -way associativity can be used to alleviate such conflicts. However, if a working set contains $p > w$ conflicts on some cache line, set associativity can only eliminate at most w of those conflicts. Our studies suggest that when conflict misses dominate, the critical factor is not a lack of associativity, but a defective line placement algorithm which fails to disperse data equitably between the available cache lines.

3.1 XOR-mapping schemes

The use of XOR-mapping schemes has been studied extensively in the context of interleaved memories [8], [15], [21], [11], [16], [10], [17] and [23] among others. In this paper we consider two types of XOR-based mapping schemes; those chosen in an *ad hoc* way based on common intuitive notions of how such schemes behave, and a scheme proposed by Rau [17] which describes a method for constructing XOR mapping schemes based on polynomial arithmetic.

The former type of XOR-mapping computes a cache index by performing a bitwise XOR of two fields of the address of the requested data. We will refer to this type of schemes as *bitwise XOR mapping*. The family of mapping functions proposed in [18] belong to this category.

In this paper we refer to Rau's scheme simply as *polynomial mapping*. Polynomial mapping can be understood by first considering address $A = \langle a_{n-1}, \dots, a_1, a_0 \rangle$ as a polynomial $A(x) = a_{n-1}x^{n-1}, \dots, a_1x^1, a_0$, the coefficients of which are in the Galois Field GF(2). The use of polynomial arithmetic, with coefficients restricted in this way, ensures that multiplication and addition of coefficients takes place modulo 2, and thus can be

implemented as logical AND and exclusive-OR respectively. The mapping from an address to an l -bit cache index is determined by the polynomial $R(x)$ defined by $A(x) = V(x)P(x) + R(x)$, where $P(x)$ is an irreducible polynomial of order l and $P(x)$ is such that $x^i \bmod P(x)$ generates all polynomials of order lower than l . The polynomials that fulfill the previous requirements are called *I-Poly* polynomials. Rau shows how the computation of $R(x)$ can be accomplished by the vector-matrix product of the address and an $n \times l$ matrix H of single-bit coefficients. In GF(2), this product is computed by a network of AND and XOR gates, and if the H -matrix is constant the AND gates can be omitted and the mapping then requires just l XOR gates with fan-in from 2 to n .

The choice of an I-poly polynomial yields properties similar to prime integer modulus functions. Whereas a prime integer modulus function would be prohibitively complex, the I-poly polynomial modulus function has very low complexity; suitable even for computing a cache index.

The use of XOR-mapping schemes requires the computation of several XOR operations to obtain the cache index. Since all the XOR can be done in parallel, the delay of this computation is just one XOR gate. The XOR gates have just two inputs for the bitwise XOR scheme and a few more for the polynomial mapping scheme. However, the computation of these XOR operations can be done at the end of the address computation stage of the pipeline. In many current microprocessors, this stage is not the critical stage of the pipeline and therefore this delay may not affect the pipeline cycle time. In addition, if some kind of carry propagate adder is used, the address computation unit computes the address bits from least-significant to most-significant. Since the XOR-mapping schemes only use some of the least-significant bits of the address, the XOR gates can operate in parallel with the computation of the most significant-bits and their delay could be completely hidden even if the address computation stage was the critical stage of the pipeline. In other situations, the addition of this small additional delay can affect the critical path, but even in these cases, a net benefit could be obtained since the reduction in miss ratio achieved by XOR-mapping schemes is very high as we will show in this paper. An accurate timing evaluation is required in these cases to consider the additional delay, which is beyond the scope of this paper.

3.2 Cache memories

This paper evaluates the performance of XOR-mapping schemes for a number of cache organizations: direct-mapped, two-way associative, victim, hash-rehash, column-associative and two-way skewed-associative. Direct-mapped and set-associative organizations [20] are the most popular in current microprocessors and we assume that the reader is familiar with them. The two-way skewed-associative cache was described in section 2. Below there is a short outline of the victim, hash-rehash and column-associative caches.

The hash-rehash cache, proposed by Agarwal *et al.* [3], consists of a conventional direct-mapped cache for which up to two tag probes may be required to find the requested data. First, the cache is accessed with the conventional modulo function, that is using l bits of the address $\langle a_{b+l-1}, a_{b+l-2}, \dots, a_b \rangle$ (2^l is the number of cache lines and 2^b is the line size). If the data is not found, the cache is probed again but with the most significant bit inverted. Thus, the second probe checks the tag for line $\langle \bar{a}_{b+l-1}, a_{b+l-2}, \dots, a_b \rangle$. In case of a second probe hit, the two lines are swapped. Otherwise, the data is brought from the next memory level and it is placed in the first-probe location, whereas the data already there is moved to the second-probe location.

The column-associative cache [4] improves the disappointing miss ratio of the hash-rehash cache by introducing a rehash bit associated with each line. This bit indicates whether the line contains rehashed data, that is, data that is reached in the second probe. When the first probe finds rehashed data, the corresponding

line is chosen for replacement. If the rehash bit is zero, then upon a first-time miss the cache is accessed again with the second function. In the case of a second-time hit, the lines are swapped. Otherwise, the data retrieved from memory is placed in the first line and the data already in that line is moved to the line accessed with the second function.

A victim cache [14] consists of a conventional direct-mapped cache with a small fully-associative buffer in the refill path to a second-level cache or main memory. On a cache miss, the line that is evicted from the direct-mapped cache is placed in the victim cache. In the case of a miss in the direct-mapped cache that hits in the victim cache, the lines accessed in both caches are swapped. In the experiments performed in this paper, we assume a victim cache with four lines.

3.3 Evaluation methodology

The results presented in this paper have been obtained through simulation of various data cache organizations using the SPEC 95 floating point benchmark suite. We focus on the floating point benchmarks because they exhibit a much higher conflict miss ratio than the integer benchmarks, and thus the XOR-mapping schemes have more potential benefits for them. Integer benchmarks will also benefit from XOR-mapping schemes although to less extent. The performance metrics used for comparison of different schemes are the total miss ratio and conflict miss ratio. Since the compared schemes only differ on the placement function, a reduction in the miss ratio will result in a reduction in the average memory access time.

The programs were compiled with the maximum optimization level and instrumented with the ATOM tool [22]. A data cache memory similar to the first-level cache of the Alpha 21164 microprocessor has been assumed: 8 Kilobytes capacity, 32 bytes per line, write-through and no write allocate. For each benchmark we have simulated the first billion (2^{30}) load operations. Because of the no write allocate feature, the performance metrics computed below refer only to load operations.

4 Performance of conventional mapping schemes

Table 1 shows the miss ratio for the following cache organizations: direct-mapped, two-way associative, four-way associative, hash-rehash, column-associative victim and two-way skewed-associative. Of these schemes, only the two-way skewed-associative cache uses an XOR-mapping scheme, as proposed by its author. For comparison, the miss ratio of a fully-associative cache is shown in the penultimate column. For each organization, the difference between its miss ratio and that of a fully-associative cache, which is shown in brackets in Table 1, represents the conflict miss ratio [12]. In fact, this difference is slightly negative in the case of 104.hydro2d and 141.apsi for some organizations, due to sub-optimality of LRU replacement in a fully-associative cache for these particular programs. Effectively the conflict miss ratio represents the target reduction in miss ratio that we hope to achieve through improved mapping schemes. The other type of misses, compulsory and capacity, will remain unchanged by the use of the XOR-mapping schemes.

From the results in Table 1, we can conclude that set associativity reduces the miss ratio, as expected, although the improvement of a two-way associative cache over a direct-mapped cache is rather low. Comparing the direct-mapped and two-way associative cache with the fully-associative cache suggests that, several benchmarks (e.g. 101.tomcatv, 102.swim, 125.turb3d, 146.wave) show significant clustering in the mapping of memory lines to cache lines under the conventional mapping scheme.

The hash-rehash cache has a miss ratio similar to that of a direct-mapped cache. Although both have similar access times, the hash-

	direct	2-way	4-way	hash-rehash	col- assoc.	victim	2-way skew	fully- assoc.
101.tomcatv	53.8 (41.3)	48.1 (36.4)	29.5 (17.0)	51.4 (39.1)	47.0 (34.5)	26.6 (14.1)	22.1 (9.6)	12.5
102.swim	56.2 (48.3)	59.1 (51.2)	57.1 (49.2)	57.6 (49.7)	53.7 (45.8)	33.7 (25.8)	15.1 (7.2)	7.9
103.su2cor	11.0 (2.1)	9.1 (0.2)	9.0 (0.1)	11.1 (2.2)	9.3 (0.4)	9.5 (0.6)	9.6 (0.7)	8.9
104.hydro2d	17.6 (0.1)	17.1 (-0.4)	17.3 (-0.2)	17.6 (0.1)	17.2 (-0.3)	17.0 (-0.5)	17.1 (-0.4)	17.5
107.mgrid	3.8 (0.3)	3.6 (0.1)	3.5 (0.0)	6.1 (2.6)	4.2 (0.7)	3.7 (0.2)	4.1 (0.6)	3.5
110.applu	7.6 (1.7)	6.4 (0.5)	6.0 (0.1)	7.8 (1.9)	6.5 (0.6)	6.9 (1.0)	6.7 (0.8)	5.9
125.turb3d	7.5 (4.7)	6.5 (3.7)	5.3 (2.5)	7.7 (4.9)	6.4 (3.6)	7.0 (4.2)	5.4 (2.6)	2.8
141.apsi	15.5 (3.0)	13.3 (0.8)	11.3 (-1.2)	18.0 (5.5)	13.4 (0.9)	10.7 (-1.8)	11.5 (-1.0)	12.5
145.fpppp	8.5 (6.8)	2.7 (1.0)	2.1 (0.4)	5.9 (4.2)	2.7 (1.0)	7.5 (5.8)	2.2 (0.5)	1.7
146.wave	31.8 (17.9)	31.7 (17.8)	23.0 (9.1)	35.4 (21.5)	30.7 (16.8)	20.1 (6.2)	16.8 (2.9)	13.9
Average	21.32 (12.61)	19.76 (11.05)	16.42 (7.71)	21.87 (13.16)	19.11 (10.40)	14.27 (5.56)	11.05 (2.34)	8.71

Table 1 Miss ratios (%) for the original schemes. The conflict miss ratio is shown in brackets.

rehash scheme requires two cache probes for some hits. Hence, the direct-mapped cache will be more effective. This poor behavior of the hash-rehash cache was also observed in [4]. The column-associative cache provides a miss ratio similar to that of a two-way associative cache. Since the former has a lower access time but requires two cache probes to satisfy some hits, the choice between these two organization should take into account the particular implementation parameters (access time and miss penalty). The victim cache removes many conflict misses and it outperforms a four-way associative cache. Finally, the two-way skewed-associative cache offers the lowest miss ratio, which is significantly lower than that of a four-way associative cache. The results for the skewed-associative cache are more positive than those observed in [18], where a miss ratio similar to a four-way associative cache was claimed, though using a different workload.

5 Bitwise XOR mapping

XOR-mapping schemes exhibit a behavior which is in some way similar to full associativity but with some restrictions. For instance, in the two-way skewed-associative cache, the set of all addresses that are mapped into the same line of bank 0 are distributed over all the lines in bank 1. Thus, it is similar to having all the lines of bank 1 as alternative locations for a given line in bank 0. However, if one considers a particular memory address, it can be placed in exactly two cache locations (one in bank 0, and the other in bank 1). Below we analyze the performance of bitwise XOR mapping schemes for the other cache organizations. The mapping functions that are evaluated are based on the family of functions proposed in [18]. Section 7 evaluates the performance of the polynomial mapping scheme proposed in [17].

5.1 Direct-mapped

To describe the bitwise XOR mapping function, let us consider a memory address $A = \langle a_{n-1}, a_{n-2}, \dots, a_0 \rangle$ composed of the following

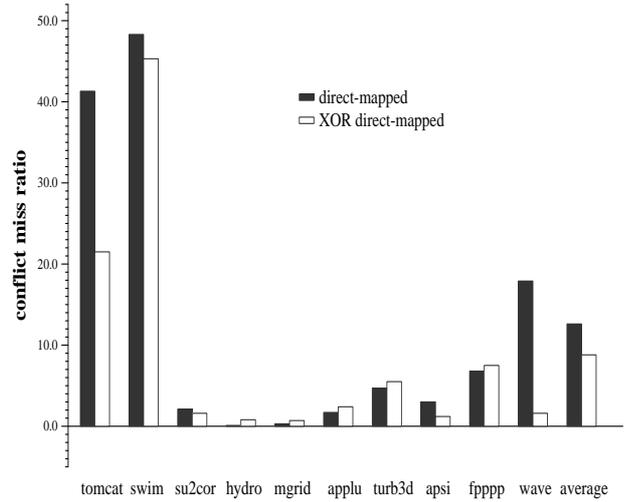


Figure 1: Conflict miss ratio (%) for a direct-mapped cache with the conventional and a bitwise XOR mapping. The average total miss ratios are 21.32% and 17.51% respectively.

fields: $A = \langle A_3, A_2, A_1, A_0 \rangle$ such that $A_0 = \langle a_{b-1}, \dots, a_0 \rangle$; $A_1 = \langle a_{l+b-1}, \dots, a_b \rangle$; $A_2 = \langle a_{2l+b-1}, \dots, a_{l+b} \rangle$; and $A_3 = \langle a_{n-1}, \dots, a_{2l+b} \rangle$. Let \oplus denote the bitwise exclusive OR. The XOR-based mapping function is defined as follows:

$$f: \{0 \dots 2^n - 1\} \rightarrow \{0 \dots 2^l - 1\}$$

$$A = \langle A_3, A_2, A_1, A_0 \rangle \rightarrow \langle A_2 \oplus A_1, A_0 \rangle$$

Figure 1 compares the conflict miss ratio of a direct mapped cache with a conventional mapping function to a direct-mapped cache with the mapping function f previously defined. It can be seen that the use of an XOR-mapping function provides a large improvement for two of the benchmarks (tomcatv and wave). These are the two benchmarks that also most benefit from a low degree of set-associativity, as can be seen from Table 1. Notice however, that five of the ten programs exhibit slightly higher miss ratios. These are all notable for their low conflict miss ratios in a conventional direct-mapped cache. We are seeing the random introduction, with low probability, of conflicts that were not originally present. On average, the direct-mapped cache with an XOR-mapping function has a total miss ratio (17.51) lower than that of a column associative cache (19.11) and almost equal to the miss ratio of a four-way associative cache (16.42).

5.2 Hash-rehash and column-associative

The mapping functions proposed for the skewed-associative cache [18] can also be used for a hash-rehash cache and a column associative cache. For these, as for the skewed-associative cache, we define two distinct mapping functions f_0 and f_1 . The first probe uses f_0 and, if required, the second probe uses f_1 . These functions are as defined in section 2, using the address decomposition $A = \langle A_3, A_2, A_1, A_0 \rangle$ defined in section 5.1, and with a binary value of $T = 10101010$.

The average total cache miss ratios for hash-rehash and column associative caches using f_0 and f_1 are 23.63% and 20.39% respectively. On average the XOR-mapping functions do not provide any improvement although they are beneficial for two benchmarks (101 and 146). The net deterioration in miss ratio is due to two reasons:

- If reference A produces a cache miss, it is placed in $f_0(A)$. If the data currently in this location corresponds to memory address B , it is moved to $f_1(A)$, or discarded. The hash-rehash cache always moves the data, whereas the column-associative cache takes this decision based on the rehash bit.

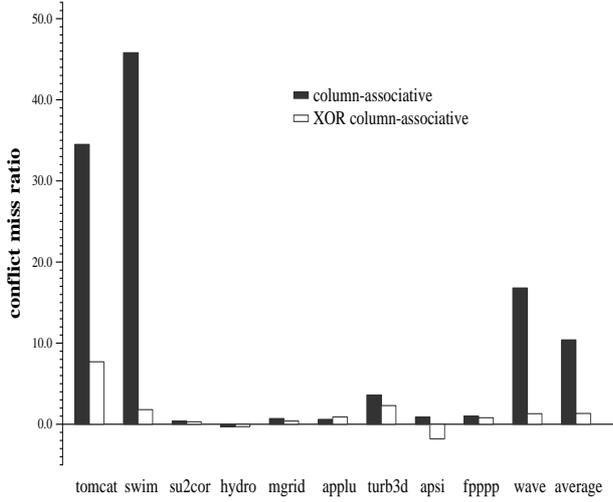


Figure 2: Conflict miss ratio (%) for the conventional column-associative cache and the new bitwise XOR-mapping. The average total miss ratios are 19.11% and 10.04% respectively.

However, it is very likely that $f_1(A) \neq f_0(B)$ and $f_1(A) \neq f_1(B)$. Consequently, the data from address B will be moved to a place where it will no longer be accessible and the next reference to B will miss (even if the data is in cache). In addition to degrading performance, this may also cause some consistency problems.

- For a given reference, it may happen that $f_0(A) = f_1(A)$. In this case, reference A does not have an alternative location and we lose the positive effect of pseudo-associativity caused by the use of two mapping functions.

5.2.1 Enhancing the hash-rehash and column-associative cache

The first problem mentioned above can be solved by inhibiting the swapping of data. Of course, that will cause a significant increase in the percentage of hits that require two probes, but it will provide us with a lower bound on the miss ratio that could be obtained. Besides, swapping may significantly increase pressure on the cache ports, and may cause performance penalties as it is not always possible to hide the swapping during idle cache cycles. For instance, in an ideal out-of-order machine with two memory ports and infinite resources, we have measured that on average two memory ports are busy during 71% of cycles, and only in the 17% of cycles are both idle [9]. An interesting alternative to swapping is to predict the most likely location of the two possible candidates for a given address. This has been extensively studied by Calder *et al.* who showed that it can be a very effective approach [6].

In order to eliminate the possibility that $f_0(A) = f_1(A)$, we propose to slightly modify the mapping functions such that they always differ in the most significant bit of the result they produce. This most significant bit will be equal to the most significant bit of A_1 for f_0 and it will be inverted for f_1 .

The proposed replacement policy is a pseudo-LRU policy inspired by the one proposed in [18]. A one-bit flag is associated with each cache line. When a hit occurs, the flag of the line holding the data is reset to 0 and the flag of the alternate location is set to 1. If a miss occurs, the new line replaces the line whose flag is higher. If both flags are equal, the line at $f_0(A)$ is replaced.

With these changes to the mapping function and replacement policy, and the elimination of swapping, the conflict miss ratios for a column-associative cache are as shown in Figure 2. The figure

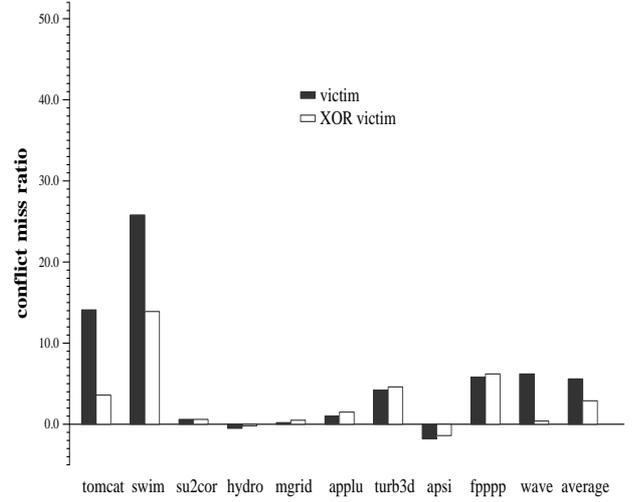


Figure 3: Conflict miss ratio (%) for the victim cache with the conventional and the bitwise XOR mapping functions. The average total miss ratios are 14.3% and 11.6% respectively.

also includes the conflict miss ratios for the conventional column-associative organization without using a XOR-mapping.

Notice that with this organization, the effect of the XOR-mapping scheme in the column-associative cache is very impressive, in particular for those programs with the highest miss ratio. The average total miss ratio of this organization (10.04) is much lower than that of a four-way associative cache (16.42) and somewhat lower than that of the skewed-associative cache (11.05). To isolate the effect of inverting one bit to obtain always two potential locations for each address, we have performed the simulations just with the XOR-functions (f_0 and f_1 as defined at the beginning of this section), without the bit inversion and we obtained an average total miss ratio of 11.17.

5.3 Victim cache

In this case, the direct-mapped part uses the XOR-mapping function defined in section 5.1. The results are shown in Figure 3. We can see that the XOR-mapping makes the average total miss ratio of the victim cache (11.6) to be very close to that of the two-way skewed-associative cache (11.05). Notice also that the XOR-mapping produces a slight increase in miss ratio for those benchmarks with very few conflict misses. The same behavior was observed for a direct-mapped cache and can be explained again by the random, but infrequent, introduction of new conflict misses.

5.4 Two-way associative

In the case of a two-way associative cache, consider an address A composed of four fields $A=(A_3, A_2, A_1, A_0)$ of $n-2l-b+2$, $l-1$, $l-1$ and b bits respectively. In this case, the XOR-based mapping function is defined as follows:

$$g: \{0 \dots 2^n - 1\} \rightarrow \{0 \dots 2^{l-1} - 1\}$$

$$A=(A_3, A_2, A_1, A_0) \rightarrow (A_2 \oplus A_1, A_0)$$

The same mapping function is used to access both banks, as in a conventional set-associative cache, and LRU replacement is used as in this case it can be implemented with low cost. The conflict miss ratios corresponding to this organization are shown in Figure 4.

The bitwise XOR mapping scheme more than halves the average total miss ratio (from 19.76 to 9.54). We can also see in Figure 4 that the mapping function has eliminated almost all the conflict misses. In average, the total miss ratio is just 1.10 times that of a fully-associative cache. For two programs the two-way XOR cache has lower miss ratio than a fully-associative cache. This is again due

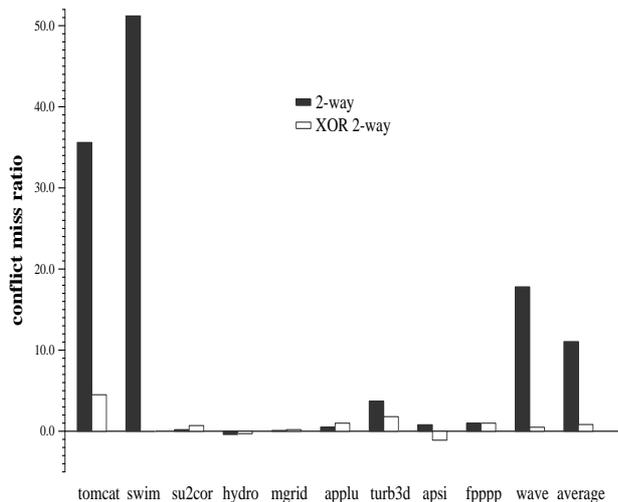


Figure 4: Conflict miss ratio (%) for the two-way associative cache with the conventional and the bitwise XOR mapping functions. The average total miss ratios are 19.76% and 9.54% respectively.

to the sub-optimality of LRU replacement in the fully-associative cache, and is a common anomaly in programs with negligible conflict misses.

When a bitwise XOR mapping is used, the average total miss ratio of the two-way associative cache (9.54) is slightly better than that of a column-associative cache (10.04) and much better than that of the skewed associative cache (11.05). This may seem to contradict the results in [18], where Seznec observed that the two-way skewed-associative cache had a lower miss ratio than a two-way associative cache with the same mapping function for both banks. The reason for this difference is twofold. Firstly, Seznec used function f_0^T described in section 2 to index the two-way associative cache. This function indexes the cache using $l - 1 + \lfloor (l - 1) / 2 \rfloor$ bits, whereas his two-way skewed-associative cache was indexed using $2l - 2$ bits. One of the most important benefits of XOR-mapping schemes is that they avoid conflicts among data structures that are accessed simultaneously with the same stride but whose initial addresses differ in a sum of powers of two. If these powers of two correspond to bits that are used by the mapping function, the conflicts may be avoided. Thus, to be fair, one should compare cache organizations that use the same number of bits as input to the mapping function. Both the two-way skewed-associative cache in Table 1 and the two-way associative cache in Figure 4 use the same number of bits. The second reason for the difference with Seznec’s results is that he used a different workload, with a much smaller working set, since his miss ratios are much lower.

The results in Figure 4 suggest that in the case of a two-way associative cache, it is more effective the use of more bits in each mapping function than having two different indexing functions.

5.5 Restricted hashing

A drawback of the XOR-mapping scheme is that it may interfere with the use of a physically tagged cache, which may be desirable for coherency reasons [12]. To remove address translation from the critical path it is common to have a virtually-indexed cache with physical address tags. This typically means that the cache is indexed using only unmapped virtual address bits. This limits the maximum number of sets and therefore, it imposes some constraints in both the cache size and the degree of associativity.

However, the XOR-mapping scheme requires the use of more bits of the address and therefore, heightens the constraint on the page size. One way to overcome this problem is to use fewer bits to compute the mapping. In the case of a skewed associative cache, it was shown that this produces a small reduction in performance [18]. We have evaluated the miss ratio of a column-associative cache using the mapping functions described in section 5.2 with the miss ratio obtained when using only the four least-significant bits of A_2 to perform the bitwise XOR with A_1 . The results showed an increase in average by a factor of 1.08, which is relatively low.

6 An affordable implementation of LRU replacement

The use of two different XOR-mapping functions creates an effect similar to full associativity, as previously discussed. This suggests that an LRU replacement policy may be expensive to implement, and has motivated previous work on pseudo-LRU replacement policies [18]. However, implementing LRU replacement in column-associative or skewed-associative caches is not as expensive as in the case of a fully-associative cache. One way to implement LRU for the caches that use two different mapping functions is to add a time stamp to each cache line. A count of memory references is maintained, and every time a cache line is accessed its time stamp is updated with the value of the reference counter. When a miss occurs, the candidate for replacement which has the lowest time stamp is chosen for replacement. In the case of a two-way skewed-associative or a column-associative cache this requires a single comparison between two integer fields.

This replacement policy produces a noticeable benefit in the performance of the column-associative cache and the two-way skewed-associative cache, as shown in Table 2, especially for benchmarks 101 and 102.

miss ratio	column-associative		2-way skewed ass.	
	ps-LRU	LRU	ps-LRU	LRU
101.tomcatv	20.2	16.4	22.1	20.0
102.swim	9.7	8.6	15.1	12.3
103.su2cor	9.2	9.0	9.6	9.1
104.hydro2d	17.2	17.1	17.1	17.1
107.mgrid	3.9	3.9	4.1	3.9
110.applu	6.8	6.4	6.7	6.3
125.turb3d	5.1	4.6	5.4	4.9
141.apsi	10.7	10.0	11.5	10.5
145.fpppp	2.5	2.5	2.2	2.2
146.wave	15.2	14.6	16.8	16.3
Average	10.04	9.31	11.05	10.24

Table 2: Miss ratios (%) for the column-associative cache and the two-way skewed-associative cache comparing pseudo-LRU with LRU replacement.

The cost associated with this LRU replacement depends on the number of bits devoted to the time-stamp. The simulations reported in Table 2 ensure that the time-stamp never overflows. A more practical scheme, that uses a small number of bits both in the counter and the time-stamp would work by shifting the counter and all the time-stamps one bit to the right whenever the reference counter overflowed. We simulated this scheme for the column-associative cache using just 8 bits for the counter and the time stamps. The results are practically identical to those obtained with an unrestricted time stamp (the average miss ratio was 9.32).

One potential criticism of our comparison between the column-associative and the skewed-associative caches is that the former

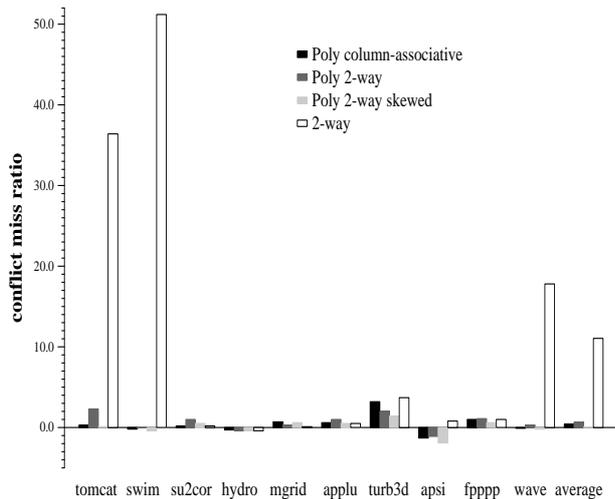


Figure 5: Conflict miss ratio for the column-associative, two-way associative and skewed-associative organizations with polynomial mapping. The conflict miss ratio of a conventional 2-way associative cache is also depicted for comparison.

uses one bit more of the address to compute the cache index. To isolate this effect we simulated the column-associative cache using $2l-2$ address bits in the mapping function (the same as the skewed-associative cache), and without bit inversion. This produced an average miss ratio of 9.36, indicating no significant difference.

7 Polynomial mapping

We have investigated the performance of the XOR-mapping scheme proposed by Rau [17], which is based on polynomial arithmetic and which will be referred to as polynomial mapping. The performance of polynomial mapping has been evaluated for the column associative, the two-way associative and the two-way skewed-associative organizations. For all of them, Table 3 compares the total miss ratios of the previous XOR mapping functions based on the bitwise XOR of two bit strings (*XOR*) with that obtained using polynomial mapping functions (*Poly*). In all cases, an LRU replacement is assumed. The miss ratio of a fully-associative cache is also shown for comparison.

To perform a fair comparison we applied the randomization scheme using the same number of bits of the original address as input to all the mapping functions; in all the cases this is 19 bits (14 without considering the bits that indicate the displacement inside the cache line). For the polynomial mapping functions, we chose the I-poly polynomials that require the fewest number of XOR entries for its implementation. We refer to a polynomial by the value obtained after substituting x by 2 (e.g., polynomial 19 is x^4+x+1). The four chosen polynomial are: $P_1=505$, $P_2=301$, $P_3=131$, $P_4=137$. For the column-associative cache, P_1 and P_2 define the mapping of the two indexing functions used by this organization. P_3 corresponds to single function utilized by the two-way associative cache. Finally, P_3 and P_4 define the two different mapping functions used by skewed-associative cache. Each mapping function requires 7 or 8 XOR gates with fan-in from 2 to 5 each.

Regarding the column-associative and the two-way associative results, we can conclude from Table 3 that the scheme based on using polynomial mapping provides a marginal advantage over the bitwise XOR scheme. However, as the former requires wider XOR gates (i.e. more inputs) the simpler XOR scheme may be preferable.

miss ratio	column-associative		2-way associative		2-way skewed assoc.		fully-assoc.
	XOR	Poly	XOR	Poly	XOR	Poly	
101.tomcatv	13.8	12.8	17.0	14.8	20.0	12.6	12.5
102.swim	8.3	7.7	7.9	7.9	12.3	7.5	7.9
103.su2cor	9.1	9.1	9.6	9.9	9.1	9.4	8.9
104.hydro2d	17.1	17.2	17.2	17.1	17.1	17.1	17.5
107.mgrid	4.0	4.2	3.7	3.8	3.9	4.1	3.5
110.applu	6.6	6.5	6.9	6.9	6.3	6.4	5.9
125.turb3d	5.5	6.0	4.6	4.8	4.9	4.2	2.8
141.apsi	10.6	11.2	11.4	11.4	10.5	10.6	12.5
145.fpppp	4.0	2.7	2.7	2.8	2.2	2.3	1.7
146.wave	14.7	13.8	14.4	14.2	16.3	13.7	13.9
Average	9.36	9.12	9.54	9.37	10.24	8.78	8.71

Table 3: Miss ratios for a column associative cache, a two-way associative cache and a two-way skewed-associative cache for the two XOR-mapping schemes: bitwise XOR (*XOR*) and polynomial mapping (*Poly*).

The marginal advantage of the polynomial mapping scheme can be explained in a number of ways. Firstly, both schemes are really quite similar; the principal advantage of polynomial mapping is the guarantee of optimal behavior on address patterns that lead to pathological conflict misses in a conventional mapping scheme. Such optimality may not be a feature of bitwise XOR schemes, but pathological cache behavior is also not a dominant feature of the SPEC95 suite. Anyway, both schemes achieve a miss ratio that is very close to that of a fully-associative cache.

On the other hand, the polynomial mapping provides a significant improvement for the skewed-associative cache. For three of the benchmarks (101, 102 and 146) this improvement is quite important. For the others, the reduction in miss ratio is very small, if any, since the miss ratio of the original mapping was already very close to that of a fully-associative cache. Overall, the skewed-associative cache using polynomial mapping and a pure LRU replacement achieves a miss ratio practically identical to that of a fully-associative cache (it is just 0.8% higher).

Figure 5 shows the conflict miss ratio for the column-associative, two-way associative and skewed-associative organizations with polynomial mapping. It can be seen that the in the three cases, practically all conflict misses have been removed.

8 Swapping in the column-associative cache

In the previous sections, the column-associative cache did not incorporate the swapping feature. As a result we can expect a lower miss ratio but a higher percentage of hits requiring two probes. We have compared the performance of the column-associative cache both with and without swapping, using a bitwise XOR mapping scheme taking $2l-2$ bits. In this case, when a reference to address A misses in cache, it is brought to $f_0(A)$. If B is the address of the data currently in that location, either it is moved to its alternative location ($f_0(B)$ or $f_1(B)$) or it is discarded if its alternative location has been used more recently. In the same way, when data is found in the second probe ($f_1(A)$) it is moved to $f_0(A)$ and the data currently in this location is moved or discarded following the same criteria as in the case of miss. In any case, data is always placed in an accessible location. We have observed that swapping increases the average total miss ratio by a factor of 1.14, but also ensures that almost all hits can be achieved with a single probe (96%).

9 Conclusions

We have analyzed the performance of XOR-based placement functions for cache memories using the SPEC 95 floating-point benchmark suite. We have shown that XOR-mapping schemes provide a very high improvement across a broad range of different cache organizations: direct-mapped, set-associative, column-associative and victim cache. We have also evaluated their effect on the hash-rehash cache and presented performance measures of the skewed-associative cache.

The main conclusion of this study is that XOR-based placement functions significantly reduce the number of conflict misses for all cache organizations. In particular, XOR-mapping combined with (pseudo) two-way associativity eliminates practically all the conflict misses, and obtains a miss ratio practically equal to that of a fully associative cache.

We have also presented a low-cost implementation of LRU replacement suitable for caches with two or more distinct mapping functions based on XOR-mapping schemes, and shown that it yields significant improvement over previously proposed pseudo-LRU replacement schemes.

Two class of placement functions have been considered. The first one is based on the bitwise exclusive OR of two bit strings. The second class is the polynomial mapping proposed in [17] in the context of interleaved memories.

For the first class of mapping functions, among the different schemes evaluated, the lowest miss ratio is achieved by the column associative cache, closely followed by the two-way set associative cache, the two-way skewed-associative cache and the victim cache. All of them achieve a miss ratio much lower than that of a conventional four-way associative cache and close to that of a fully-associative cache. For example, a two-way associative cache achieves an average miss ratio that is just 1.09 times that of a fully-associative cache. Similarly, a column-associative cache can achieve a miss ratio between 1.07 and 1.23 times that of a fully-associative cache, depending on whether swapping is implemented. For comparison, a conventional direct-mapped cache has a miss ratio that is 2.45 times that of a fully-associative cache.

Regarding polynomial mapping, we have shown that it provides a marginal advantage over the simpler bitwise XOR schemes for the two-way associative and column-associative organizations. However, for the skewed-associative cache it achieves a significant reduction in miss ratio. Combining the effects of a LRU replacement and polynomial mapping, the miss ratio of the two-way skewed associative cache is reduced from 1.27 to 1.01 times that of a fully associative cache.

Comparing the three most effective organizations, i.e., skewed-associative, column-associative and set-associative, we can see that all achieve a very similar miss ratio. Each one may be preferable for different reasons: a skewed-associative has the lowest miss ratio, the column-associative has the lowest hit time and the set-associative requires less hardware to implement a LRU replacement.

In overall, we can conclude that XOR-based placement functions are an extremely powerful technique for eliminating conflict misses.

Acknowledgments

This work has been supported by the Spanish Ministry of Education (grants CICYT TIC-429/95 and Acción Integrada Hispano-Británica 202B); the British Council (grant 1016); and the UK EPSRC (grant K19723).

We would like to thank the anonymous referees for their constructive comments.

References

- [1] Amdhal Corp., *470V/6 Machine Reference Manual*, 1976
- [2] A. Agarwal, *Analysis of Cache Performance for Operating Systems and Multiprogramming*, Kluwer Academic Publishers, 1989, pp. 120-122.
- [3] A. Agarwal, J. Hennessy and M. Horowitz, "Cache Performance of Operating Systems and Multiprogramming", *ACM Trans. on Comp. Systems*, 6, Nov. 1988, pp. 393-431.
- [4] A. Agarwal and S.D. Pudar, "Column-Associative Caches: A Technique for Reducing the Miss Rate of Direct-Mapped Caches", in *Proc. Int. Symp. on Computer Architecture*, 1993, pp. 179-190.
- [5] T. Asprey *et al.*, "Performance Features of the PA7100 Microprocessor", *IEEE Micro*, 13(3), June 1993, pp. 22-35.
- [6] B. Calder, D. Grunwald and J. Emer, "Predictive Sequential Associative Caches", in *Proc. Int. Symp. on High Performance Computer Architecture*, 1996, pp. 244-253.
- [7] D.A. Fotland *et al.*, "Hardware Design of the First HP Precision Architecture Computer", *Hewlett-Packard Journal*, 38(3), March 1987, pp. 4-17.
- [8] J.M. Frailong, W. Jalby and J. Lenfant, "XOR-Schemes: A Flexible Data Organization in Parallel Memories", in *Proc. Int. Conf. on Parallel Processing*, 1985, pp. 276-283.
- [9] J. González and A. González, "Identifying Contributing Factors to ILP", in *Proc. Euromicro 96*, 1996
- [10] D.T. Harper III, "Reducing Memory Contention in Shared Memory Multiprocessors", in *Proc. Int. Symp. on Computer Architecture*, 1991, pp. 66-73.
- [11] D. Harper III and D. Linebarger, "A Dynamic Storage Scheme for Conflict-Free Vector Access", in *Proc. Int. Symp. on Computer Architecture*, 1989, pp. 72-77.
- [12] J.L. Hennessy and D. Patterson, *Computer Architecture: A Quantitative Approach*, Morgan Kaufmann Publiss., 1996.
- [13] IBM, *3033 Processor Complex, Theory of Operation/Diagrams Manual-Processor Storage Control Function*, vol. 4, IBM, Poughkeepsie, N.Y., 1978
- [14] N. P. Jouppi, "Improving Direct-Mapped Cache Performance by the Addition of a Small Fully-Associative Cache and Prefetch Buffers", in *Proc. Int. Symp. on Computer Architecture*, 1990, pp. 364-373.
- [15] A. Norton and E. Melton, "A Class of Boolean Linear Transformations for Conflict-free Power-of-two Stride Access", in *Proc. Int. Conf. on Parallel Processing*, 1987, pp. 247-254.
- [16] B.R. Rau, M.S. Schlansker and D.W.L. Yen, "The Cydra 5 Stride-Insensitive Memory System", in *Proc. Int. Conf. on Parallel Processing*, 1989, pp. 242-246.
- [17] B.R. Rau, "Pseudo-Randomly Interleaved Memories", in *Proc. Int. Symp. on Computer Architecture*, 1991, pp. 74-83.
- [18] A. Seznec, "A Case for Two-way Skewed-associative Caches", in *Proc. Int. Symp. on Computer Architecture*, 1993, pp. 169-178.
- [19] A. Seznec and F. Bodin, "Skewed-associative Caches", in *Proc. Int. Conf. on Parallel Architectures and Languages (PARLE)*, 1993, pp. 305-316.
- [20] A. J. Smith, "Cache Memories", *ACM Computing Surveys*, 14(4), Sept. 1982, pp. 473-530.
- [21] G. S. Sohi, *Logical Data Skewing Schemes for Interleaved Memories in Vector Processors*, Computer Science Technical Report #753, Univ. of Wisconsin-Madison, Sept. 1988.
- [22] A. Srivastava and A. Eustace, "ATOM: A System for Building Customized Program Analysis Tools", in *Proc. SIGPLAN Conf. on Programming Language Design and Implementation*, 1994.
- [23] M. Valero *et al.*, "Increasing the Number of Strides for Conflict-free Vector Access", in *Proc. Int. Symp. on Computer Architecture*, 1992, pp. 372-381