

# Near-Optimal Loop Tiling by means of Cache Miss Equations and Genetic Algorithms

Jaume Abella, Antonio González, Josep Llosa  
Computer Architecture Department  
Universitat Politècnica de Catalunya  
Barcelona (Spain)  
{jabella, antonio, josepll}@ac.upc.es

Xavier Vera  
Institutionen för Datateknik  
Mälardalens Högskola  
Västerås (Sweden)  
xavier.vera@mdh.se

## Abstract

*The effectiveness of the memory hierarchy is critical for the performance of current processors. The performance of the memory hierarchy can be improved by means of program transformations such as loop tiling, which is a code transformation targeted to reduce capacity misses. This paper presents a novel systematic approach to perform near-optimal loop tiling based on an accurate data locality analysis (Cache Miss Equations) and a powerful technique to search the solution space that is based on a genetic algorithm. The results show that this approach can remove practically all capacity misses for all considered benchmarks. The reduction of replacement misses results in a decrease of the miss ratio that can be as significant as a factor of 7 for the matrix multiply kernel.*

## 1 Introduction

Memory performance is critical for the performance of current computers. Memory is organized hierarchically in such a way that the upper levels are smaller and faster. The uppermost level typically has a very short latency (e.g. 1-2 cycles) but the latency of the lower levels may be a few orders of magnitude longer (e.g. main memory latency may be around 100 cycles). Thus, techniques to keep as much data as possible in the uppermost levels are key to performance.

In addition to the hardware organization, it is well known that the performance of the memory hierarchy is very sensitive to the particular memory reference patterns of each program. The reference patterns of a given program can be changed by means of transformations that do not alter the semantics of the program. These program transformations can modify the order in which some computations are performed or can simply change the data layout. *Loop Tiling*

is an example of the former family of techniques. Loop tiling [1, 2] is a technique based on a combination of strip-mining and loop interchange.

Loop tiling has a significant potential to remove cache misses. Loop tiling can remove most capacity misses by restructuring the loop and changing the order in which statements are executed. However, finding the optimal loop tiling for a given program is a very complex task, since the options are almost unlimited and exploring all of them is infeasible. For very simple programs, the programmer intuition may help but in general, a systematic approach that can be integrated into a compiler and can deal with any type of program is desirable. This systematic approach requires the support of a locality analysis method in order to assess the performance of different alternatives.

In this paper, we propose an automatic approach to perform loop tiling in numeric codes. It is based on a very accurate technique to analyze the locality of a program that is known as Cache Miss Equations (CMEs) and a genetic algorithm in order to search the solution space. The proposed genetic algorithm converges very fast, and although it does not guarantee that the optimal solution is found, we show that after loop tiling, the replacement miss ratio of the evaluated benchmarks is almost negligible. Therefore, for these programs the results are near-optimal. The rest of this paper is organized as follows. Section 2 overviews the locality analysis approach. Section 3 presents the loop tiling technique and its performance is evaluated in section 4. Section 5 outlines some related work and section 6 summarizes the main conclusions of this work.

## 2 Memory Locality Analysis

To effectively transform a code in order to optimize memory performance, an effective memory locality analysis is required in order to assess the different alternatives. In this section, we describe the locality analysis technique

```

parameter (N)
REAL a(N,N), b(N,N), c(N,N)
do i = 1, N
  do j = 1, N
    do k = 1, N
      a(i,j) = a(i,j) + b(i,k) * c(k,j)
    enddo
  enddo
enddo

```

**Figure 1. Matrix multiply algorithm**

used in this work to perform loop tiling.

In particular, we use Cache Miss Equations (CMEs) [3] to represent the cache behavior. Cache Miss Equations are a very accurate analytical model of the cache memory. They describe the cache behavior by means of diophantine equations, which allows us to use mathematical techniques to compute the locality of each memory reference. Unfortunately a direct solution of these equations is computationally intractable due to its NP nature. Statistical methods and some techniques based on polyhedra theory have been proposed to solve the equations in a reasonable amount of time [4, 5].

## 2.1 Overview of Cache Miss Equations

Cache Miss Equations are a set of equations<sup>1</sup> that represent all the potential cache misses for the references in a loop nest. They describe the precise relationship among the iteration space, array sizes, base addresses, and the cache parameters for a loop nest. This section presents an overview of the CMEs. For more details about CMEs see [3, 6].

In order to generate CMEs, the *reuse vectors* [7] of all the references in a loop nest must be generated. *Reuse vectors* provide information about the potential reuses in the entire iteration space. Figure 1 shows the *matrix multiply* kernel. For instance,  $\vec{r} = (0, 0, 1)$  is a reuse vector for reference  $c(k,j)$ , because the data accessed by this reference in a given iteration is potentially reused one iteration later (both are potentially mapped into the same cache line). In order to determine if these potential reuses are realized, CMEs are generated and studied (it is not necessary to solve the equations to find realized reuses).

For every reuse vector of a reference two types of CMEs are generated:

- **Compulsory equations.** Compulsory equations represent the first time a memory line is brought into the cache.

<sup>1</sup>The term equation is loosely used to refer to a set of simultaneous equalities and inequalities.

- **Replacement equations.** Given a reference, replacement equations represent the interferences with any other reference. For each pair of references ( $R_A$  and  $R_B$ ), the following expression gives the condition that determines whether they are mapped onto the same cache set:

$$Cache\_Set(\vec{i})_{R_A} = Cache\_Set(\vec{j})_{R_B} \\ \vec{j} \in \mathcal{I}$$

where  $\mathcal{I}$  represents the iteration points between  $\vec{i}$  (the current one) and the iteration point from which  $R_A$  reuses. This condition is expanded into a set of equations for each reuse vector.

## 2.2 Finding Cache Misses from CMEs

Deciding whether a reference causes a miss or a hit for a given iteration point is equivalent to deciding whether this iteration point belongs to the polyhedra defined by the CMEs. The points inside each CMEs polyhedron represent the potential cache misses (the number of points is the number of potential cache misses). This leads us to consider several ways for computing them:

- **Solver.** Given a reference  $R$  with  $m$  reuse vectors and  $n_k$  equations for the  $k^{th}$  reuse vector, the polyhedron that contains all the iteration points that result in a miss is [3]:

$$Set\_Misses = \cap_{k=1}^m \cup_{j=1}^{n_k} Solution\_Set\_Equation_j$$

This approach implies counting the number of points inside the union of convex polyhedra. This requires counting the points (which is an NP problem for a single polyhedron) in an exponential number of polyhedra, making this problem infeasible due to its huge computing time.

- **Traversing the iteration space.** Given a reference, all the iteration points can be tested independently [6]. In order to know if an iteration point  $\vec{i}_0$  results in a miss we need to know when it fulfills the CMEs. This problem is equivalent to finding out whether, after substituting the iteration point in the CMEs, the resulting polyhedron is non-empty. This is still an NP problem, but only a linear number of polyhedra must be analyzed for each iteration point. However, the problem is still infeasible except for tiny iteration spaces.

Moreover, in a  $k$ -way set associative cache, there are  $k$  cache lines in every set, so  $k$  distinct contentions are needed before a cache miss occur. Therefore, the first method can only be applied to direct-mapped caches whereas the second method works for both direct-mapped and set-associative organizations.

## 2.3 A Fast and Accurate Implementation to Solve CMEs

In this section we describe a fast and accurate approach to estimate the solution to the CMEs. Our approach builds upon the second method to solve the CMEs (traversing the iteration space). A key property of CMEs is that each point of the iteration space and each memory reference can be studied independently of the others.

We have developed some techniques that exploit the special characteristics of the CMEs polyhedra [4] in order to speed-up the process of counting points in polyhedra:

- **Counting Compulsory Polyhedra.** When an iteration point is substituted, Compulsory Equations result in polyhedra with either 0 or 1 variable. A polyhedron with 0 variables consists of a set of inequality relations between integer values. The iteration point that has been substituted is a potential miss if the inequalities hold. On the other hand, since a polyhedron with 1 variable represents an interval, the iteration point  $\vec{t}_0$  results in a potential miss when there exist integer values in it.
- **Counting Replacement Polyhedra.** Due to the particular form of these polyhedra, the number of integer solutions can be computed in a more efficient way than in general polyhedra. We have developed a method to detect when they are empty that is based on counting the number of integer points inside them [8]. In order to compute it, the domains<sup>2</sup> of the different variables involved in its definition are calculated. This can be done by means of the vertices of the polyhedron, but computing them is an NP problem. We have developed specific techniques for replacement polyhedra that compute the domains of the variables in a polynomial time. A detailed description of these techniques can be found in [8].

The use of these techniques results in an average speed-up of 20 over a method based on identifying the vertices of the polyhedra. However, this important speed-up is still insufficient to solve CMEs for huge iteration spaces in a reasonable amount of time.

To further reduce the computation cost, we use sampling techniques to study a subset of the iteration space instead of the whole iteration space [5] [9]. The subset of points is selected using *Simple Random Sampling* [10].

We model the number of misses of each reference using a Discrete Random Variable. This random variable follows a Binomial distribution that models phenomena consisting of  $n$  different and independent experiments that fol-

<sup>2</sup>We define the real domain of a variable  $x$  in a polyhedron  $P$  as the range of real values it takes inside  $P$ .

low a Bernoulli distribution. We can use statistical techniques [11] in order to compute the parameters that describe this random variable. The approach to obtain an approximation of the miss ratio is to evaluate the behavior of a subset of the population (sample) obtaining the empirical value of the parameters that describe the sample and to infer these values to the population.

The size of the sample is set according to the required width of the confidence interval and the desired confidence. We found experimentally that a confidence interval of width 0.1 and a 90% confidence is enough to obtain accurate miss ratios with a very small computing time (only 164 points of the iteration space must be explored). In this way, the miss ratio is computed as an interval of width 0.1 and the actual miss ratio belongs to this interval with a probability of 90%. The central point of this interval can be used as an estimation of the actual miss ratio.

## 2.4 CME for Multiple Convex Regions

CME are defined for iteration spaces that are single convex regions, but after tiling  $n$  dimensions the iteration space is the union of  $2^n$  convex regions.

Figure 2 shows how the iteration space of a one-dimensional loop becomes a two-convex region iteration space after tiling. The shaded regions correspond to the different convex regions before and after tiling.

There are different ways to solve this problem. The easiest option would be to use only one convex region that approximates the actual non-convex region. This convex region can be the smallest parallelepiped that includes all other convex regions (see figure 2 (c)) or alternatively, the region which does not include the last iteration of every tiled loop where the tile size is not a divisor of the upper bound (see figure 2 (d)). Both options have drawbacks. The first option includes in the convex regions points outside the iteration space, while the second option does not include points belonging to the iteration space.

Because of this, we have decided to implement a more accurate solution. The CME implementation has been modified to deal with multiple convex regions by defining the equations for every convex region and solving for every analyzed point the equations corresponding to the convex region in which the point is contained. Let  $n$  be the number of convex regions of a loop after tiling. Every compulsory equation should be defined for each convex region, so the number of compulsory equations is increased by a factor of  $n$ . For each reuse vector, we have to generate a set of replacement equations for each convex region. In addition, we have to generate a set of equations for every pair of convex regions that reflect the potential reuse between different regions. Due to this, the number of replacement equations is increased by a factor of  $n^2$ .

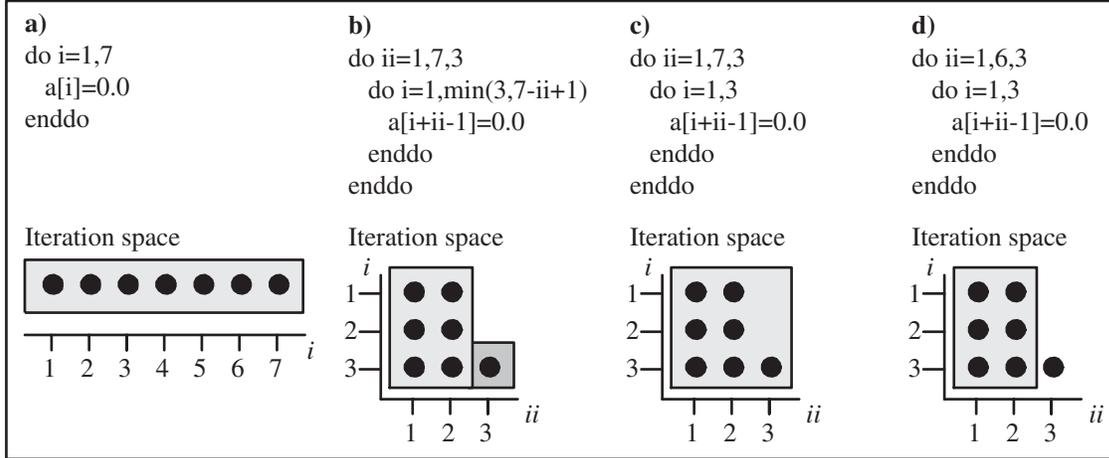


Figure 2. Example of iteration space: (a) before tiling, (b) after tiling, (c) bigger region, (d) smaller region

<p>(a)</p> <pre>do i<sub>1</sub>=1, U<sub>1</sub>   do i<sub>2</sub>=1, U<sub>2</sub>     A(i<sub>2</sub>,i<sub>1</sub>) = B(i<sub>1</sub>,i<sub>2</sub>)   enddo enddo</pre>	<p>(b)</p> <pre>do ii<sub>1</sub>=1, U<sub>1</sub>, T<sub>1</sub>   do ii<sub>2</sub>=1, U<sub>2</sub>, T<sub>2</sub>     do i<sub>1</sub>=ii<sub>1</sub>, min(ii<sub>1</sub>+T<sub>1</sub>-1, U<sub>1</sub>)       do i<sub>2</sub>=ii<sub>2</sub>, min(ii<sub>2</sub>+T<sub>2</sub>-1, U<sub>2</sub>)         A(i<sub>2</sub>,i<sub>1</sub>) = B(i<sub>1</sub>,i<sub>2</sub>)       enddo     enddo   enddo enddo</pre>
---	---

Figure 3. 2D matrix transposition: (a) before tiling, (b) after tiling

### 3 Loop Tiling

In this section we present our proposal to perform loop tiling. Tiling [12, 1] is a transformation which combines strip-mining [1] with loop interchange to form small tiles [1] of loop iterations in order to increase the data locality. In this section  $T_i$  and  $U_i$  stand for the tile size and upper bound respectively of loop  $i$  in the original loop nest. Figure 3 shows an example of a loop before (a) and after (b) loop tiling.

#### 3.1 Model

The target of our work is to obtain the values of the variables  $T_i$  that minimize the number of misses. Note that loop

tiling changes only the order in which the original iteration space is traversed, so the number of compulsory misses before and after tiling remains constant. Because of this, our work focuses on minimizing the number of replacement misses, which include both capacity and conflict misses.

Let  $f$  be the function that represents the number of replacement misses for each possible value of the tile size variables:

$$f \mapsto \#Replacement\ Misses \quad (1)$$

$$f : \underbrace{[1, U_1]}_{T_1} \times \dots \times \underbrace{[1, U_k]}_{T_k} \mapsto \mathbb{Z}$$

Our problem can be expressed as follows:

$$MIN \quad f(T_1, \dots, T_k)$$

$$1 \leq T_i \leq U_i$$

$$i = 1 \dots k$$

where  $f$  is called the *objective function*.

The objective function consists of the CMEs generated in a parameterized way, where the parameters are the tile sizes.

Since  $f$  is a pseudo-polynomial function [13], the relationship between loop tiling and the number of misses is nonlinear.  $T_i$  can take only integer values, thus, our problem can be seen as a nonlinear integer optimization (NLP) one.

Many researchers have studied NLP [14, 15]. A well studied case is the one where the constraints are linear (named *linearly constrained optimization*). A special case

is when the objective function is entirely linear; this is called Linear Programming (LP). Algorithms to solve both real (e.g simplex [16]) and integer functions (e.g branch & bound scheme [17]) can be found in the literature.

One of the challenges in NLP is that some problems exhibit local minima and search algorithms can be stuck at them. Algorithms that propose to overcome this problem are named *Global Optimization*. Real functions have been studied deeply [18, 19, 15]. Unfortunately, integer functions are hard to optimize. There are some studies based on  $\{0,1\}$  valued integer functions [20], but in general, this is a hard and time-consuming problem. Hence, the use of heuristics is necessary. Tabu search [21] obtains promising theoretical results, but only partial implementations have been reported so far. On the other hand, simulated annealing [22] and genetic algorithms [23, 24] have been used for years with very good results.

Our proposal is based on the use of a genetic algorithm to optimize function  $f$ .

### 3.2 Genetic Algorithm

Algorithms for function optimization are generally limited to convex regular functions. However, there are lots of functions that are not continuous, non differentiable or multi-modal. It is common to solve this problem by means of stochastic sampling. Whereas traditional search techniques use characteristics of the problem to determine the next sampling point (e.g Gradient), stochastic methods use non-deterministic decision rules [25].

Genetic Algorithms (GAs) are a particular type of stochastic methods that have been used to solve hard problems with objective functions that do not meet the properties required by traditional methods [23]. These algorithms search in the solution space of a function simulating the Nature-based process of evolution, that is, the survival of the fittest. Usually, the fittest individuals tend to reproduce more than the inferior individuals, and they survive to the next generation propagating the best genes.

GAs simulate the evolution of a population. Figure 4 shows the simplest GA. It starts from a random generated population, and it makes the population evolve by means of basic genetic operators (selection, mutation and crossover) [23] applied to individuals of the current population, to produce an improved next generation.

### 3.3 Genetic Algorithm Parameters

The use of GAs requires the determination of the following issues: chromosome representation, selection function, genetic operators, the creation of the initial population and the termination criteria.

```

ALGORITHM:
Supply a population  $P_0$ 
 $i=1$ 
while (not finish)
   $P_i=Selection(P_{i-1})$ 
   $P_i=Reproduce(P_i)$ 
   $i=i+1$ 
end

```

Figure 4. Simple Genetic Algorithm

In our context, each individual represents a particular tiling solution, which is identified by the particular values of the tile size variables  $T_i$ . Individuals are made up of a set of chromosomes, each one being associated to a tile size variable. Each chromosome is made up of a sequence of genes, each gene being represented by a digit of a certain alphabet. We have experimentally observed that using the alphabet  $\{00, 01, 10, 11\}$  produces good results.

The function to transform the chromosome values into tile sizes is not the identity function. Tile size  $T_i$  can take any value in the range  $[1 \dots U_i]$ . On the other hand a chromosome is represented by a sequence of genes encoded in a binary representation. Thus, each chromosome will be represented by a value in the range  $[0 \dots 2^k - 1]$  where  $k$  is  $\lceil \log_2 U_i \rceil$ . If  $k$  is an odd number,  $k$  is increased in 1 because of the alphabet we have used to represent genes. Thus, there are more values in the representation range for a chromosome than possible tile size values. Therefore, we need a function to map values  $[0 \dots 2^k - 1]$  into the range  $[1 \dots U_i]$ .

Let  $g$  be the function that represents the tile size for each possible value of a chromosome:

$$g : [0 \dots 2^k - 1] \rightarrow [1 \dots U_i]$$

where

$$k = \lceil \log_2 U_i \rceil \quad (+1 \text{ if odd})$$

$$x \in [0 \dots 2^k - 1]$$

$$g(x) = \lfloor \frac{x * (U_i - 1)}{2^k - 1} \rfloor + 1 \quad (2)$$

It can be deduced that every possible tile size has at least one representation.

**Example.** Let us assume a two-dimensional nested loop where the upper bounds of the loops are 10 and 100 respectively. Thus,  $\lceil \log_2 10 \rceil = 4$  and  $\lceil \log_2 100 \rceil = 7$ , so we have  $k_1 = 4$  and  $k_2 = 8$  respectively. Therefore, the first chromosome is represented by 2 genes, and the second one by 4 genes. For instance, the value 12 (1100) and 74 (01001010) correspond to the tile sizes 8 ( $g_1(12)=8$ ) and 29 ( $g_2(74)=29$ ) respectively, and are represented by the following genes:

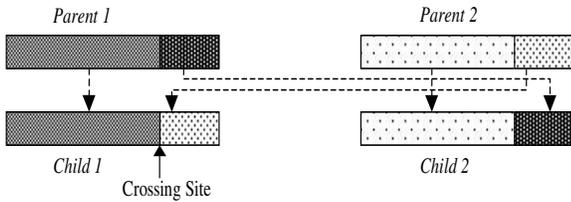
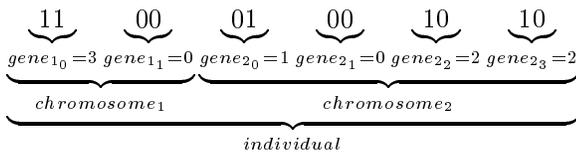


Figure 5. Schematic of simple crossover



Genetic operators provide the basic search mechanism by creating new solutions based on the solutions that exist. The *selection* of individuals to produce successive generations plays an extremely important role. A common selection approach assigns probability of selection to each individual depending on its fitness. Individuals with higher fitness have a higher probability of contributing one or more offsprings to the next generation. We have adopted one of the selection schemes that gives better results, which is known as *remainder stochastic selection without replacement* [23]. Let us call  $N$  the size of the population (number of individuals). This selection scheme consists in choosing  $N$  individuals from the  $N$  individuals of the previous generation. In this selection process a given individual can be chosen more than once. The chosen individuals are grouped forming pairs and crossover is applied to each pair with a given probability. In the case they do not crossover, both individuals are added to the new population (see Figure 5).

*Crossover* takes two individuals and produces two new individuals merging the genetic material in a random point (named cross site). *Mutation* is applied after crossover to each individual with a given probability. Mutation changes one individual to produce a new one by flipping some of its genes. Both crossover probability and mutation probability have to be determined empirically, and are related to the size of the population. Figure 6 shows an example of an iteration of the genetic algorithm.

The GA must be provided with an initial population (see Figure 4) that is created randomly. GAs moves from generation to generation, and the usual termination criterion is the number of generations, although other criteria can be used [23].

Our experiments show that an initial population of size equal to 30, with crossover probability of 0.9 and a mutation probability of 0.001, gives near-optimal results in most

Description	Indiv.	Indiv.	Indiv.	Indiv.
Generation $i$	00 00 00	01 01 01	10 10 10	11 11 11
Selection	00 00 00	11 11 11	01 01 01	00 00 00
Crossover?	Yes		No	
Crossover	00 00 11	11 11 00	01 01 01	00 00 00
Mutation	00 00 11	11 01 00	01 01 01	00 00 11
Generation $i + 1$	00 00 11	11 01 00	01 01 01	00 00 11

Figure 6. Example of a genetic algorithm iteration

cases after 15 generations. In the rest of the cases, near-optimal results are obtained after a number of generations between 15 and 25. Figure 7 shows the algorithm to decide the number of generations required before the optimal tile search stops, where *converge()* is a function to decide when the population is homogeneous enough. In our case we consider that a population converges when the best individual has a difference of replacement misses smaller than 2% with respect to the population average of its generation. We have observed in the evaluated loops that this convergence criterion is only achieved if the population is close to the optimal.

Finally, note that CME have an exponential cost with respect to the number of dimensions of the loop nest, and loop tiling is a transformation that doubles the number of dimensions of the loop nest and increases the number of equations. In spite of this, due to the efficient implementation of CMEs, the required 450 evaluations (15 iterations of the GA  $\times$  30 individuals) for each loop nest can still be solved in a reasonable time. In our case, every loop nest took between 15 minutes and 4 hours on a SUN sparc Ultra-60 workstation. This compilation time can be assumed for those codes which performance is crucial, like scientific and embedded processor applications.

## 4 Performance Evaluation

This section evaluates the proposed loop tiling approach. Examples of its use will be given, as well as its accuracy.

### 4.1 Experimental Framework

CMEs have been implemented for fortran codes through the Polaris Compiler [26] and the Ictineo library [27]. These libraries allow us to obtain all the compile-time information needed to generate the equations.

The evaluation of CMEs has been implemented in C++ following the techniques outlined in section 2.3 and using our own polyhedra representation [4].

Due to CMEs restrictions, only perfectly nested loops in

```

ALGORITHM:
finish := false
iters := 0
while (not finish)
  if (iters < 15)
    iters = iters + 1
    create next generation
  else if (iters >= 15 and iters < 25)
    if (not converge())
      iters = iters + 1
      create next generation
    else finish := true
  endif
  else finish := true
  endif
endwhile

```

**Figure 7. Genetic Algorithm used to perform loop tiling**

which the array subscript expressions are affine functions of the induction variables are analyzed [3].

The loop nests considered are some kernels from different programs (NAS<sup>3</sup>, BIHAR<sup>4</sup>, LIVERMORE) and some frequently used kernels (see Table 1). These loops have been chosen because they exhibit high number of capacity misses. Results for different cache architectures are reported.

## 4.2 Examples of some Kernels

First of all we show the effectiveness of the loop tiling technique by means of some well-known kernels (see Table 2). We have evaluated their miss ratio for a 8KB direct-mapped cache with 32-byte lines.

Table 2 shows the results. Column 2 shows the problem size. Columns 3 and 4 show the total and replacement miss ratio before loop tiling respectively, whereas columns 5 and 6 show the total and replacement miss ratio after tiling respectively. We can see that after tiling the replacement miss ratio is near zero for all kernels. This indicates that loop tiling has removed almost all replacement misses. The remaining replacement misses probably are conflict misses which could be removed by means of techniques like padding.

## 4.3 Performance Evaluation for some Kernels

In this section we present the replacement miss ratio of different direct-mapped caches (8KB and 32KB) before and after applying loop tiling for all the benchmarks of the table 1. Results are shown in Figures 8 and 9, where the

<sup>3</sup>Numerical Aerospace Simulation Facility

<sup>4</sup>Biharmonic Partial differential equations solver

Kernel	Program	Nested loops	Description
T2D		2	2D Matrix transposition
T3DJIK		3	3D Matrix transposition $a[k,j,i] = b[j,i,k]$
T3DIKJ		3	3D Matrix transposition $a[k,j,i] = b[i,k,j]$
JACOBI3D		3	Partial differential equations solver
MATMUL		3	Matrix by vector multiplication
MM	LIVER.	3	Matrix multiplication
ADI	LIVER.	2	2D ADI integration
ADD	NAS	4	Addition of update to a matrix
BTRIX	NAS	3	Block Tri-diagonal solver. Backward block sweep
VPENTA1	NAS	2	Invert 3 pentadiagonals simultaneously. Loop 1
VPENTA2	NAS	2	Invert 3 pentadiagonals simultaneously. Loop 2
DPSSB	BIHAR	3	unnormalized inverse of a forward transform of a complex periodic sequence
DPSSF	BIHAR	3	forward transform of a complex periodic sequence
DRADBG1	BIHAR	3	backward transform of a real coefficient array. Loop 1
DRADBG2	BIHAR	3	backward transform of a real coefficient array. Loop 2
DRADFG1	BIHAR	3	forward transform of a real periodic sequence. Loop 1
DRADFG2	BIHAR	3	forward transform of a real periodic sequence. Loop 2

**Table 1. Evaluated kernels**

number following the program name corresponds to the problem size.

It can be seen that for most programs practically all replacement misses are removed after tiling, which implies that near-optimal solutions have been found for these programs. However, for some kernels (ADD, BTRIX, VPENTA1 and VPENTA2) the replacement miss ratio obtained after tiling is still quite high for all cache sizes. For some others (ADI\_1000 and ADI\_2000) replacement misses after tiling are only significant for a 8KB cache. We have analyzed these cases carefully and we have observed that most of the remaining replacement misses are due to conflicts and they cannot be removed by loop tiling. For these programs we have investigated the combination of padding and tiling techniques. Padding parameters are obtained in a similar way to tiling ones. They are introduced in the CMEs and a GA is used to find near-optimal solutions. Details can be found in [28].

Table 3 shows the replacement miss ratios obtained for those kernels where the replacement miss ratio is still high after loop tiling. Column 2 shows the original replacement miss ratio; column 3 lists the replacement miss ratio after

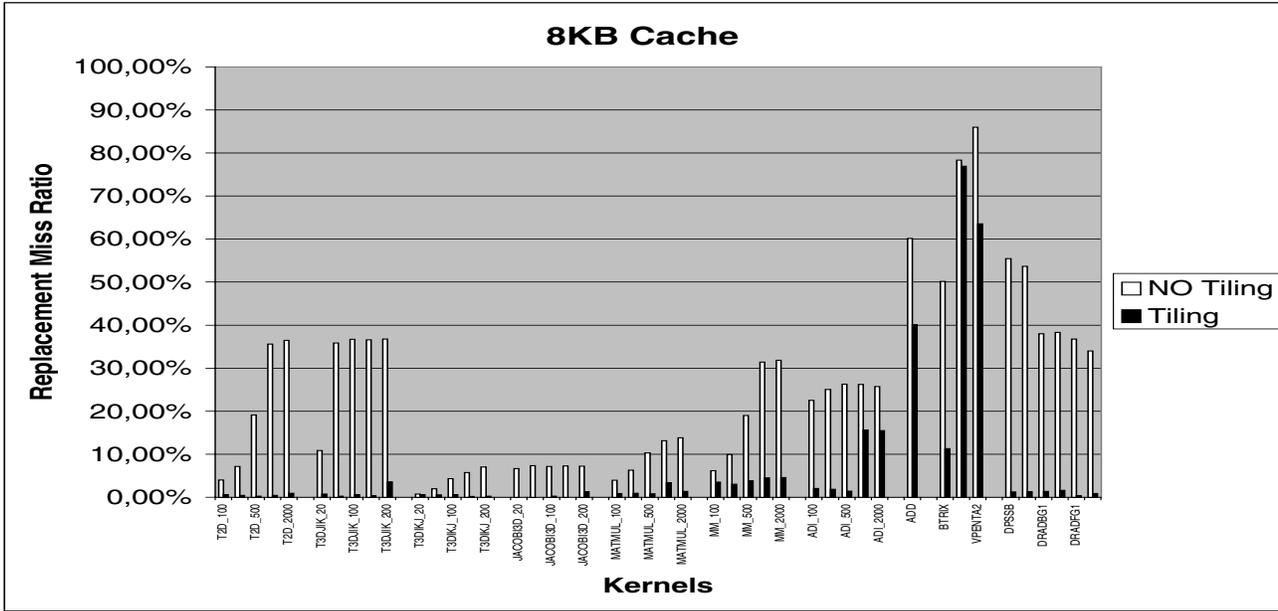


Figure 8. Replacement miss ratio before and after loop tiling for a 8KB cache.

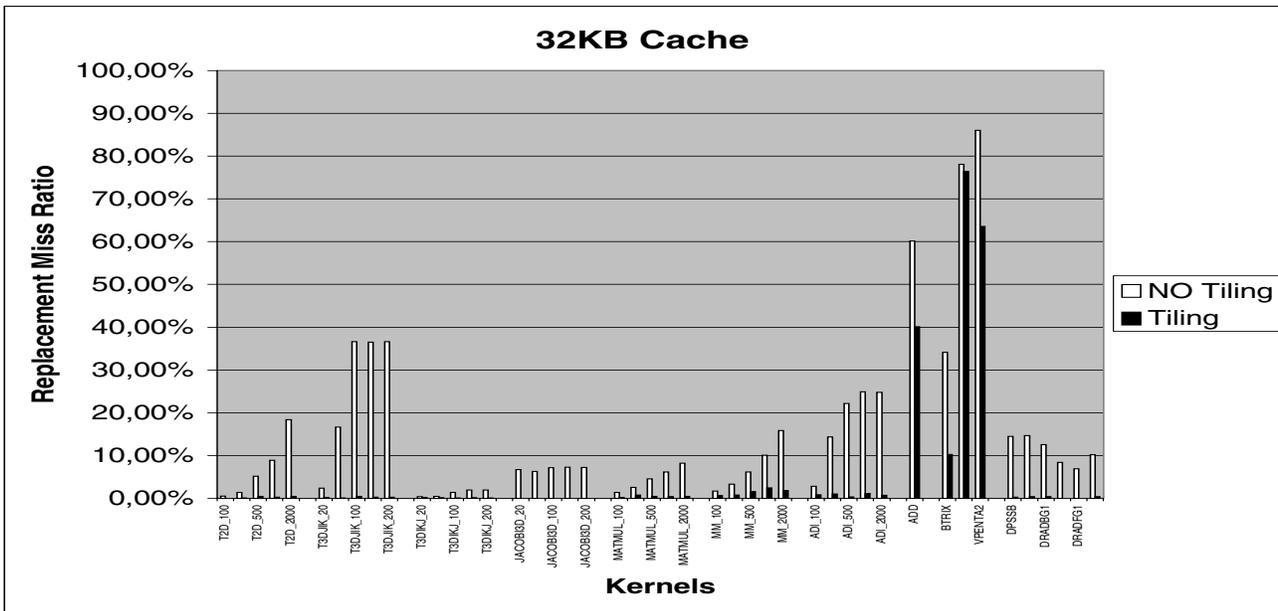


Figure 9. Replacement miss ratio before and after loop tiling for a 32KB cache.

padding; and column 4 shows the replacement miss ratio after padding and tiling are applied sequentially in this order. In the future we plan to study the application of padding and tiling techniques in a single step, trying to find the padding and tiling parameters at the same time. This can in general produce better results than optimizing each part separately.

The replacement miss ratios after padding and tiling are practically null for all benchmarks.

Table 4 shows the percentage of kernels (not considering those that appear in table 3), which have a replacement miss ratio lower than 1%, 2% and 5% respectively after tiling. For all these kernels and all cache sizes the replacement

Kernel	Prob size	No Tiling		Tiling	
		Total	Repl.	Total	Repl.
T2D	N=2000	63.3%	36.4%	27.7%	0.9%
T3DJIK	N=200	63.4%	36.7%	30.2%	3.6%
T3DIKJ	N=200	34.6%	7.0%	27.9%	0.3%
JACOBI3D	N=200	25.6%	7.2%	19.8%	1.3%

**Table 2. Miss ratio for some evaluated kernels (8KB direct-mapped cache, 32B lines)**

Kernel	8KB		
	Original	Padding	Padding + tiling
ADD	60.2%	59.8%	0.5%
BTRIX	50.1%	0.2%	0.2%
VPENTA1	78.3%	52.4%	0.0%
VPENTA2	86.0%	11.9%	0.0%
ADL1000	26.2%	12.3%	4.1%
ADL2000	25.7%	12.4%	3.4%

Kernel	32KB		
	Original	Padding	Padding + tiling
ADD	60.2%	59.8%	0.0%
BTRIX	34.1%	0.0%	0.0%
VPENTA1	78.1%	32.9%	0.0%
VPENTA2	86.0%	11.3%	0.0%

**Table 3. Miss ratio for some evaluated kernels before padding and tiling, after padding and after padding and tiling**

Cache sizes	Replacement miss ratio		
	<1%	<2%	<5%
8KB	56.4%	79.5%	100.0%
32KB	90.2%	97.6%	100.0%

**Table 4. Replacement miss ratios after tiling of all kernels excepting those in table 3**

miss ratios are lower than 5%.

Our technique is compared against the optimal solution (counting replacement misses) but not against other techniques in the literature. Due to the different limitations of these techniques they cannot be compared with the same benchmarks and same platform on an equal basis.

## 5 Related Work

Caches are an essential part of processors for reducing memory latency and increase memory bandwidth. Some

programs have loops which work with large working sets that exceed the cache capacity and result in a high number of capacity misses. These misses can be reduced by means of computation-reordering transformations such as loop interchange, loop distribution, loop skewing and loop tiling among others [1, 2].

Loop tiling [12] is an effective optimizing transformation to reduce the number of capacity misses of programs, especially for dense matrix computations. However, the success of loop tiling depends on the tile size and shape selection. Many algorithms have been provided to find suitable approximations for this selection.

Ghosh, Martonosi and Malik proposed the use of CME [3] to select the tile size [29] but they did not propose a general algorithm to do it. Their technique consists on maximizing the tile size for every self-interference equation, obtaining a tile that has no replacement misses for the given equation. They do not give details about how to combine the different tile sizes obtained for every self-interference equation. Their tiling algorithm is not applied to cross-interference equations.

Coleman and McKinley presented a technique that tries to maximize the tile size such that it fits in cache and at the same time it reduces cross-interference misses. Their technique to reduce cross-interference misses is based on computing the worst case: maximum number of expected cross-interference misses. To do this, their algorithm estimates the footprints of array references [30].

Rivera and Tseng proposed a tile size selection algorithm [31] for programs that compute values using neighboring array elements in a fixed stencil pattern. In order to avoid cross-interference misses they use different techniques based on copying tiles, using a subset of the cache, padding and applying the same technique as Coleman and McKinley.

Sarkar and Megiddo presented a constant-time algorithm to obtain near-optimal tile sizes for two-dimensional nested loops [32] taking into account self and cross-interference misses. The algorithm is based on an approximated memory cost model and an analytical model to estimate the memory cost of a loop nest. The algorithm is extended in order to deal with three-dimensional nested loops using an iterative search over the first tile size and applying their algorithm over the two inner loops.

Our proposal differs and improves these previous approaches in the fact that it is a technique based on a precise model to represent cache behavior, that searches the solution space for optimal tile sizes, for any type of reference pattern that corresponds to affine references and for loop nests of any dimension. It always produces a tiling scheme that reduces capacity misses and usually removes practically all replacement misses.

## 6 Conclusions

Cache memory performance is critical for the efficient execution of numerical applications. Loop tiling is a program transformation that reduce capacity misses. In this work, we have proposed the combination of genetic algorithms and a very accurate model of the cache behavior in order to perform near-optimal loop tiling. The cache model is based on a very fast solver of the Cache Miss Equations. The proposed technique can deal with any perfectly nested loop.

The evaluation of the proposed technique shows that the resulting tiling significantly reduces the miss ratio. For instance, for a 8KB direct-mapped cache, we can reduce the replacement miss ratio of the 3D matrix transposition (N=100) from 36.7% to 0.6% and the replacement miss ratio of the Dpssb kernel from 55.5% to 1.25%. We have shown that the proposed loop tiling technique practically removes all capacity misses for all the loops that have been analyzed.

## References

- [1] Susan L.Graham David F.Bacon and Olivier J.Sharp. Compiler transformations for high-performance computing. Technical report, University of California, 1994.
- [2] S.Carr K.S.Kinley and C.W.Tseng. Improving data locality with loop transformations. In *ACM Transactions on Programming Languages and Systems*, 1996.
- [3] Somnath Ghosh, Margaret Martonosi, and Sharad Malik. Cache miss equations: an analytical representation of cache misses. In *ICS97*, 1997.
- [4] Nerina Bermudo, Xavier Vera, Antonio González, and Josep Llosa. An efficient solver for cache miss equations. In *IEEE International Symposium on Performance Analysis of Systems and Software*, 2000.
- [5] Xavier Vera, Josep Llosa, Antonio González, and Carlos Ciuraneta. A fast implementation of cache miss equations. In *8th International Workshop on Compilers for Parallel Computers (CPC)*, January 2000.
- [6] Somnath Ghosh, Margaret Martonosi, and Sharad Malik. Precise miss analysis for program transformations with caches of arbitrary associativity. In *ASPLOS98*, 1998.
- [7] Michael E. Wolf and Monica S. Lam. A data locality optimizing algorithm. In *ACM SIGPLAN91*, 1991.
- [8] Nerina Bermudo, Xavier Vera, Antonio González, and Josep Llosa. Optimizing cache miss equations polyhedra. In *4th Workshop on Interaction between Compilers and Computer Architecture (Interact)*, 2000.
- [9] Somnath Ghosh, Margaret Martonosi, and Sharad Malik. Cache miss equations: a compiler framework for analyzing and tuning memory behavior. In *ACM Transactions on Programming Languages and Systems*, November 1999.
- [10] Moore; McCabe. *Introduction to the Practice of Statistics*. Freeman & Co, 1989.
- [11] M.H. DeGroot. *Probability and statistics*. Addison-Wesley, 1998.
- [12] R.Schreiber and Jack Dongarra. Automatic blocking of nested loops. Technical report, RIACS, 1990.
- [13] Philippe Clauss. Counting solutions to linear and nonlinear constraints through ehrhart polynomials: applications to analyze and transform scientific programs. In *ICS96*, 1996.
- [14] Bazaraa, Shetty, and Sherali. *Nonlinear programming: theory and applications*. Wiley, 1994.
- [15] Gill, Murray, and Wright. *Practical optimization*. Academic Press, 1981.
- [16] S.G. Nash and A. Sofer. *Linear and nonlinear programming*. McGraw Hill, 1996.
- [17] G. Sierksma. *Linear and integer programming: theory and practice*. M. Dekker, 1996.
- [18] Torn and Zilinskas. *Global optimization*. Springer-Verlag, 1989.
- [19] Host, Pardalos, and Thoai. *Introduction to global optimization*. Kluwer, 1995.
- [20] Hansen, Jaumard, and MATHON. Constrained nonlinear 0-1 programming. *ORSA Journal on Computing*, 1995.
- [21] Glover and Laguna. *Tabu search*. Kluwer, 1997.
- [22] Kirkpatrick, Gelatt, and Vecchi. Optimization by simulated annealing. *Science* 220, 1983.
- [23] D.E. Goldberg. *Genetic algorithms in search, optimizations and machine learning*. Addison-Wesley, 1989.
- [24] J. Holland. *Adaptation in natural and artificial systems*. The University of Michigan Press, Ann Arbor, 1975.
- [25] Y. Ermoliev and R.J.-B. Wets. *Numerical Techniques for Stochastic Optimization*. Springer-Verlag, 1988.
- [26] David Padua et al. *Polaris developer's document*, 1994.
- [27] Eduard Ayguadé et al. A uniform internal representation for high-level and instruction-level transformations. Technical Report UPC-DAC-95-02, Universitat Politècnica de Catalunya, 1995.
- [28] Xavier Vera, Antonio González, and Josep Llosa. Near-optimal padding for removing conflict misses. Technical Report UPC-DAC-2000-71, Universitat Politècnica de Catalunya, November 2000.
- [29] Somnath Ghosh, Margaret Martonosi, and Sharad Malik. Cache miss equations: a compiler framework for analyzing and tuning memory behavior. In *ACM Transactions*, 1998.
- [30] S.Coleman and K.S.McKinley. Tile size selection using cache organization and data layout. In *PLDI*, 1995.
- [31] Gabriel Rivera and Chau-Wen Tseng. Tiling optimizations for 3d scientific computations. In *SC'00*, 2000.
- [32] Vivek Sarkar and Nimrod Megiddo. An analytical model for loop tiling and its solution. In *ISPASS'00*, 2000.