# Reducing Memory Requirements of Nested Loops for Embedded Systems

J. Ramanujam\*

Jinpyo Hong\*

Mahmut Kandemir<sup>†</sup>

A. Narayan\*

# Abstract

Most embedded systems have limited amount of memory. In contrast, the memory requirements of code (in particular loops) running on embedded systems is significant. This paper addresses the problem of estimating the amount of memory needed for transfers of data in embedded systems. The problem of estimating the region associated with a statement or the set of elements referenced by a statement during the execution of the entire set of nested loops is analyzed. A quantitative analysis of the number of elements referenced is presented; exact expressions for uniformly generated references and a close upper and lower bound for non-uniformly generated references are derived. In addition to presenting an algorithm that computes the total memory required, we discuss the effect of transformations on the lifetimes of array variables, i.e., the time between the first and last accesses to a given array location. A detailed analysis on the effect of unimodular transformations on data locality including the calculation of the maximum window size is discussed. The term maximum window size is introduced and quantitative expressions are derived to compute the window size. The smaller the value of the maximum window size, the higher the amount of data locality in the loop.

#### 1 Introduction

An important characteristic of embedded systems is that the hardware can be customized according to the needs of a single or a small group of applications. An example of such customization is parameterized memory/cache modules whose several topological parameters (e.g., total capacity, block size, associativity) can be set depending on the data access pattern of the application at hand. It many cases, it is most beneficial to use the smallest amount of data memory that satisfies the target performance level [20]. Employing a data memory space which is larger than needed has several negative consequences. First, per access energy consumption of a memory module increases with its size [2]. Second, large memory modules tend to incur large delays, thereby increasing the data access latency. Third, large memories by definition occupy more chip space. Consequently, significant savings in energy/area/delay might be possible be being more careful on selecting a memory size.

Unfortunately, selecting the minimum memory size (without impacting performance) is not always easy. This is because data declarations in many codes are decided based on high-level representation of the algorithm being coded not based on actual memory requirements. The important point here is that not every data item (declared data location) is needed throughout the execution of the program. That is, at any given point during the execution, typically, only a portion of the declared storage space (e.g., array size) is actually needed. Therefore, the total data memory size can be reduced

Copyright 2001 ACM 1-58113-297-2/01/0006....\$5.00.

by determining the *maximum* number of data items that are live at any point during the course of execution. Two complementary steps to achieve this objective is (i) estimating the memory consumption of a given code, and (ii) reducing the memory consumption through access pattern transformations.

The problem of estimating the minimum amount of memory was recently addressed by Zhao and Malik [20]. In this paper, we present a technique that (i) quickly and accurately estimates the number of distinct array accesses and the minimum amount of memory in nested loops, and (ii) reduces this number through loop-level transformations. Nested loops are of particular importance as many embedded codes from image and video processing domains manipulate large arrays (of signals) using several nested loops. In most cases, the number of distinct accesses is much smaller than the size of the array(s) in question and the size of the loop iteration space. This is due to the repeated accesses to the same memory location in the course of execution of the loop nest. The proposed technique identifies this reuse of memory locations, and takes advantage of it in estimating the memory consumption as well as in reducing it.

The main abstraction that our technique manipulates is that of data dependence and re-use [19]. Since many compilers that target array-dominated codes maintain some sort of data dependence information, implementing our estimation and optimization strategy involves only a small additional overhead. Our experimental results obtained using a set of seven codes show that the proposed techniques are very accurate, and are capable of reducing the memory consumption significantly through high-level optimizations.

The rest of this paper is organized as follows. Section 2 presents a brief background. Section 3 presents our techniques for estimating the number of distinct references to arrays accessed in nested loops. In Section 4, we present a loop transformation technique that minimizes the maximum amount of memory required. Section 5 presents experimental results on seven benchmarks. Related work is discussed in Section 6, and Section 7 concludes with a summary.

#### 2 Background

DSP programs mainly consist of perfectly nested loops (loop nests in which every statement is inside the innermost loop) that access single and multi-dimensional arrays [19, 20]. Therefore, we will limit our discussion to these cases. In our framework, each execution of an n-level nested loop is represented using an iteration vector  $\vec{I} = (i_1, i_2, \dots, i_n)$ , where  $i_j$  corresponds to  $j^{th}$  loop from the outermost position. We assume that the array subscript expressions and loop bounds are affine functions of enclosing loop indices and loop-independent variables [19]. Each reference to an d-dimensional array U is represented by an access (or data reference) matrix  $A_D$  and an offset vector  $\vec{b}$  such that  $A_D \vec{l} + \vec{b}$  is the element accessed by a specific iteration  $\overline{I}$  [19]. The access matrix is a  $d \times n$  matrix. An array element which is referenced (read or written) more than once in a loop nest constitutes a reuse. The reuse count depends on the number of references to the array in the loop (r), the relative values of the loop nest levels and the dimensionality (d) of the array, and also the loop limits. If iterations  $\vec{i}$  and  $\vec{j}$  access the same location, we say that the *reuse vector* is  $\vec{i} - \vec{i}$ . The level of a reuse vector is the index of the first non-zero in it. Consider the following loop nest:

<sup>\*</sup>Louisiana State University, Baton Rouge, LA 70803, USA. (jxr@ee.lsu.edu)

<sup>&</sup>lt;sup>†</sup>Pennsylvania State University, State College, PA 16801. (kandemir@cse.psu.edu) Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

DAC 2001, June 18-22, 2001, Las Vegas, Nevada, USA.

for 
$$i = 1$$
 to N do  
for  $j = 1$  to N do  
for  $k = 1$  to N do  
 $\cdots U[i, k-3] \cdots$ 

The iteration vector is  $(i, j, k)^T$  (note that we write a column vector as transpose of the corresponding row vector), , the data access matrix for array U is  $\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$ , and the offset vector is  $\begin{pmatrix} 0 \\ -3 \end{pmatrix}$ .

For an array where its reference matrix is a non-square matrix (that is, if the dimension d of the array and the loop nesting n is the same) the number of times an element is referenced is at most r where r is the number of references to the array in a loop. Therefore the *reuse* for a data element is at most r - 1. For an array whose dimension is one less than the loop nest, the reuse of an element is along the direction of the null space of the access matrix and the amount of reuse due to an element depends on the loop bounds. Arrays whose number of dimensions is one less than the depth of the nested loop enclosing them are very common in DSP applications [2].

# 2.1 Data Dependences and Loop Transformations

Dependence (and reuse) analysis is critical to the success of optimizing compilers [18, 19]. We deal with sets of perfectly nested loops, whose upper and lower bounds are all linear, enclosing a loop body with affine array references. That is, each subscript of an array variable index expression must be an affine expression over the scalar integer variables of the program. We assume familiarity with definitions of the types of dependences [19].

Dependences arise between two iterations  $\vec{I}$  and  $\vec{J}$  when they both access the same memory location and one of them writes to the location [19]. Let  $\vec{I}$  execute before  $\vec{J}$  in sequential execution; the vector  $\vec{d_i} = \vec{J} - \vec{I}$  is referred to the *dependence vector* [19]. This forces sequentiality in execution. The level of a dependence vector is the index of the first non-zero element in it [19]. Let  $\vec{I} = (I_1, \ldots, I_n)$  and  $\vec{J} = (J_1, \ldots, J_n)$  be two iterations of a nested loop such that  $\vec{I} \prec \vec{J}$  (read  $\vec{I}$  precedes or executes before  $\vec{J}$  in sequential execution) and there is a dependence of constant distance  $\vec{d_i}$  between them. Applying a linear transformation T to the iteration space (nest vector) also changes the dependence matrix since  $T(\vec{J}) - T(\vec{I}) = T(\vec{J} - \vec{I}) = T \vec{d_i}$ . All dependence vectors are positive vectors, i.e., the first non-zero component should be positive. We do not include *loop-independent* dependences which are zero vectors. A transformation is legal if the resulting dependence vectors are still positive vectors [19]. Thus, the linear algebraic view leads to a simpler notion of the legality of a transformation. For an *n*-nested sequential loop, the  $n \times n$  identity matrix  $(I_n)$  denotes sequential execution order. Any unimodular transformation can be realized through reversal, interchange and skewing [17].

#### 2.2 Distinct references

The number of distinct references  $(A_d)$  can be found using dependences in the loop as shown in Figure 1. The *n*-dimensional cube (in the case of 2-nested loop, this is a square) formed by the dependence vectors as shown in Figure 1 represents the reused area (the shaded area) in the iteration space. Consider the following examples:

Example 1(a):	for $i = 1$ to 10 do
	for $j = 1$ to 10 do
	$\cdots A[i, j] \cdots A[i, j] a i + 2]$
	$\cdots A[i-5, j+2] \cdots$
Example 1(b):	for $i = 1$ to 10 do
	for $j = 1$ to 10 do
	$\cdots A[2 * i + 3 * j] \cdots$



Dependence (3, -2)

Figure 1: Iteration space for a 2-nested loop.

In both Example 1(a) and 1(b), the dependence vector is (3, -2). Note that in the first example, the dimensionality of the array is the same as that of the loop nest level, the number of references (r) is 2 and the reuse count is at most 1. (The number of times an element of the array is referenced is at most 2).

In the second example, the dimensionality is less than the loop nest level, the number of references (r) is 1, and the maximum reuse count for an element is  $\lceil 10/3 \rceil = 4$ . The total reuse (i.e., the area of the shaded region) is the same in both the examples which is  $(10-3) \times (10-2) = 56$ . Let the dependence vector be  $(d_1, d_2)$ . In general, the signs of  $d_1$  and  $d_2$  do not affect the amount of reuse. In a nested loop of size  $N_1 \times N_2$ , the amount of reuse is given by  $(N_1 - |d_1|) \times (N_2 - |d_2|)$ . We consider the cases where the dimension of the array accessed within the loop is the same as the nest level and where the dimensionality is less than the loop nesting level. These cases are commonly found in DSP codes [2].

## 2.3 Uniformly generated references and maximum window size

We assume that all the references to an array are *uniformly generated* [9, 5]. Uniformly generated references are those, for which the access matrices are the same but the offset vectors are different, i.e., the subscript functions of the different references differ only in the constants. An example of a loop with a uniformly generated references is shown below:

for 
$$i = 1$$
 to  $N_1$  do  
for  $j = 1$  to  $N_2$  do  
 $X[2i+3j+2] = Y[i+j]$   
 $Y[i+j+1] = X[2i+3j+3]$ 

Here, the two references to X are of the form 2i + 3j + constant

and both references to array Y are of the form i + j + constant. We use the notion of a reference window of an array in loop nest (which is different form the notion of the reference window of a dependence as used by [5, 9]) that allows us to deal with each distinct array as a whole and not on a per-reference-pair-to-the-array basis.

The amount of memory required is a function of the number of variables which will be accessed again in future. We now introduce a notion that is useful in this context. The *reference window*  $W_X(\vec{I})$  (where  $\vec{I} = (I_1, \ldots, I_n)$  is an iteration of the *n*-nested loop) is the set of all elements of array X that are referenced by any of the statements in all iterations  $\vec{J_1} \leq \vec{I}$  (read  $\vec{J_1}$  precedes in sequential execution or is the same as  $\vec{I}$ ) that are also referenced in some (later) iteration  $\vec{J_2}$  such that  $\vec{J_2} \succ \vec{I}$  (read  $\vec{J_2}$  follows  $\vec{I}$ ). This allows us to precisely the define those iterations which need a specific value

in local memory. The size of the window  $W_X(\vec{I})$  is the number of elements in that window. The maximum window size (MWS) is given by  $\max_{\vec{I}} |W_X(\vec{I})|$  and is defined over the entire iteration space. In the case of multiple arrays  $X_1, \ldots, X_K$ , the maximum reference window size is:  $\max_{\vec{I}} \sum_{j=1}^{K} |W_{X_j}(\vec{I})|$ 

Note that the reference window is a dynamic entity, whose shape and size change with execution. For nested loops with uniformly generated references, the maximum window size (MWS) is a function of the loop limits. The smaller the value of MWS, the higher the amount of data locality in the loop nest for the array. For simplicity of exposition, we assume that there are multiple uniformly generated references to a single array in a loop nest. The results derived here easily generalize to multiple arrays and higher levels of nesting.

## 3 Estimating the number of distinct accesses in nested loops

#### **3.1** Loops with Array Dimension d = Nesting n

With just one reference to each array in such a nest, the number of distinct accesses equals the total number of iterations. Therefore, only the case where there are multiple references to the same array. For example, in relaxation codes, this is common.

In general for r references in a loop where the array dimension is the same as the loop nesting level there are a total of  $\frac{r(r-1)}{2}$  dependences. Note that there is at least one node in the dependence graph which is a sink to the dependence vectors from each of the remaining r-1 nodes. In other words there exists a statement with r-1 direction vectors directed from each of the remaining statements. The r-1 dependences due to all the other references to this reference gives the amount of reuse. Consider a two-level nested loop in which there are r uniformly generated references. Let the dependences on one reference due to all other references be

$$\left(\begin{array}{cccc} d_{11} & d_{21} & \cdots & d_{r-1,1} \\ d_{12} & d_{22} & \cdots & d_{r-1,2} \end{array}\right)$$

The amount of reuse for that array is: reuse =  $\sum_{i=1}^{r-1} (N_1 - |d_{i1}|) (N_2 - |d_{i2}|)$  and the number of distinct elements is given by  $A_d = N_1 \times$  $N_2 \times r$  – reuse. Consider the following loop (in Example 2) where there are two uniformly generated references to the array A and the access matrix is non-singular.

**Example 2:** for i = 1 to  $N_1$ for j = 1 to  $N_2$  $S_1$ :

 $S_2$ :

Here there is a dependence (1, -2) from statement  $S_1$  to statement  $S_2$ . This dependence is used to calculate the amount of reuse for each element. The amount of reuse is  $(N_1 - 1)(N_2 - 2)$ , and the number of distinct accesses to the array A in the above loop is  $A_d = N_1 \times N_2 \times 2$  - reuse.

Example 3:	for $i = 1$ to 10
	for $j = 1$ to 10
$S_1$ :	$\cdots A[i,j] \cdots$
$S_2$ :	$\cdots A[i-1,j] \cdots$
$S_3$ :	$\cdots A[i, j-1] \cdots$
$S_4$ :	$\cdots A[i-1,j-1] \cdots$

The dependences from statement  $S_1$  to all other statements are (1,0), (0,1), (1,1). The amount of reuse is calculated as reuse = (10-1)(10-0) + (10-0)(10-1) + (10-1)(10-1) =90 + 90 + 81 = 261, and the the number of distinct accesses is:  $A_d = 10 \times 10 \times 4 - 261 = 139$ . Thus, we see that for cases where the loop nesting level is the same as the dimension of the array accesses in the loop there is only one dependence vector between a

pair of statements and the maximum reuse for a particular element is at most r-1. In other words, there are a maximum of r references to an array element.

#### **3.2** Loops with Array Dimension d = n - 1

**Single Reference** Now consider the case where the dimension of the array is at least one less than the loop nest. If d = n - 1then there is reuse along the direction of the null space vector of the access matrix.

#### Example 4

for i = 1 to 20 do for j = 1 to 10 do  $\cdots A[2i+5j+1]\cdots$ 

Here the reuse vector is (5, -2) which is the same as the dependence vector for the loop. We now look at the n dimensional cube formed by the dependence vector (in this case, a square) on the iteration space which represents the reused elements of the array. Note that all elements within the square formed by the vector is a sink to a direction vector which is a reused element by definition. Therefore, for the above example where there is a single statement, we can obtain the figure for the number of data elements reused in the array as:

reuse = 
$$(N_1 - d_{11})(N_2 - |d_{21}|) = (20 - 5)(10 - 2) = 120$$
,

and the number of distinct accesses to the array is

$$A_d = N_1 \times N_2 - \text{reuse} = 20 \times 10 - 120 = 80.$$

Now consider the case of a 2-dimensional array accessed in a three nested loop.

## Example 5:

for i = 1 to 10 do for j = 1 to 20 do for k = 1 to 30 do  $\cdots A[3i+k, j+k] \cdots$ 

Here the reuse vector is (1, 3, -3); the reuse is calculated as:

reuse = 
$$(10 - 1)(20 - 3)(30 - 3) = 4131$$
,

and the number of distinct accesses is

$$A_d = 10 \times 20 \times 30 - 4131 = 1869$$

Multiple References The case of multiple references is not discussed in this paper for lack of space.

It is important to note that our techniques is exact for uniformly generated references.

Non-uniformly Generated References The distinct elements in loops where there are non-uniformly generated references are more complex to compute. The dependence vectors for these loops are distance vectors and thus it is not possible to represent the exact reuse using the dependence vectors. Consider the following example:

## Example 6:

for i = 1 to 20 do for j = 1 to 20 do

 $S_1: \quad \cdots \quad A[3i+7j-10] \cdots \\ S_2: \quad \cdots \quad A[4i-3j+60] \cdots$ 

Here the dependences are distance vectors and an exact dependence cannot be obtained. We give lower and an upper bound on the number of distinct accesses on the array A. We have the upper and lower bounds on both the functions  $f_1 = 3i + 7j - 10$  and  $f_2 = 4i - 3j + 60$ . We have  $LB_1 \le f_1 \le UB_1, LB_2 \le f_2 \le$ 

 $UB_2$ ,  $LB_1 = 0$ ,  $LB_2 = 4$ ,  $UB_1 = 190$ ,  $UB_2 = 137$ . Therefore the smallest lower bound  $LB_{min} = 0$ , and the largest upperbound  $UB_{max} = 190$ .

The upper bound on the number of references = 190 - 0 + 1 = 191. The lower bound on the number of references = 191 - (3 - 1)(7 - 1) - (3 - 1)(7 - 1) = 179. The actual number of references is 181. So a close bound on the number of distinct accesses to the array can be obtained by the above algorithm.

## 4 Minimizing the maximum window size using transformations

Consider the following example which is a minor variant of the example from [5]:

# Example 7:

for i = 1 to 20 do

for j = 1 to 30 do

 $\cdots X[2i-3j]\cdots$ 

Eisenbeis et al. [5] mention that the cost of the window (the same as MWS) for this loop is 89. They use only two transformations: loop interchange and reversal. On applying interchange, the MWS reduces to 41. On reversal applied to the original loop, the cost becomes 86 while reversing the interchanged loop reduces the cost to 36. Using the technique presented here, the cost or MWS for this loop can be reduced to 1, i.e., all iterations accessing any element of the array X can be made consecutive iterations of an inner loop. The only dependence in this example is the vector (3, 2). We use the following legal transformation, T:

$$T = \left( \begin{array}{cc} 2 & -3 \\ 1 & -1 \end{array} \right)$$

Even though the technique in [14] can be used to derive this transformation, there are situations where the techniques presented here *improves* locality while that in [14] does not improve locality. Consider the loop shown in the next example.

#### Example 8:

for i = 1 to 25 do

for j = 1 to 10 do X[2i + 5j + 1] = X[2i + 5j + 5]

The distance vectors for this loop are: (3, -2), (2, 0), (5, -2);(3, -2) is the flow dependence, (2, 0) is an anti-dependence and (5, -2) is the output dependence vector. These are the only direct dependences. Li and Pingali use transformation matrices whose first row is either (2, 5) or (-2, -5). Any transformation that uses (2, 5) as its first row is illegal because of the distance vector (3, -2); the first component of (3, -2) after the transformation is  $((2, 5) \cdot (3, -2)^T = -4 \text{ is } < 0)$ . Similarly any transformation that uses (-2, -5) as its first row is illegal due to the distance vector (2, 0) since  $((-2, -5) \cdot (2, 0)^T = -4 \text{ is } < 0)$ . The maximum window size is 50. Li and Pingali's technique will not find any partial transformation that can be completed to a legal transformation. Where as, by applying techniques presented in the following sections, we can apply the legal transformation, T:

$$T = \left(\begin{array}{cc} 2 & 3\\ 1 & 1 \end{array}\right)$$

Applying T reduces the maximum window size to 21. A combination of reversal and interchange does not change the maximum window size from 50.

## 4.1 Effect of Transformations on Locality

Consider a nested loop with r uniformly generated references to an array X of the form:  $\lambda_1 i + \lambda_2 j + c_k$  (k = 1, ..., r) as shown

below:

#### **Example 9:**

for i = 1 to  $N_1$  do for j = 1 to  $N_2$  do  $\cdots X[\lambda_1 i + \lambda_2 j + c_1] \cdots$  $\cdots X[\lambda_1 i + \lambda_2 j + c_2] \cdots$  $\cdots$ 

 $\cdots X[\lambda_1 i + \lambda_2 j + c_r] \cdots$ We need to compute the effect of a legal unimodular transformation, T:

$$T = \left( \begin{array}{cc} a & b \\ c & d \end{array} \right)$$

on the maximum window size . In addition to legality, we require that the loop nest be *tileable* [10, 18]; this permits us to use block transfers, which are very useful to minimize the number of off-chip accesses. The optimum transformation thus satisfies two conditions:

- 1. legality condition for tiling
- 2. minimizes the maximum window size

We do not show the detailed derivation here. The maximum window size (MWS) is a function of the maximum inner loop span or *maxspan*, which is the maximum trip count of the inner loop (difference between the upper and lower limits of the inner loop) over all outer loop iterations [5].

$$MWS = \max span \times \Delta \times (\lambda_2 a - \lambda_1 b)$$
(1)

where  $\Delta$  is the determinant of the transformation matrix. The simplified expression derived for the maximum window size is:

$$MWS = \begin{cases} \left( \left| \frac{N_1 - 1}{b} \right| + 1 \right) |\lambda_2 a - \lambda_1 b| & \text{if } a - b \ge aN_1 - bN_2 \\ \left( \left| \frac{N_2 - 1}{a} \right| + 1 \right) |\lambda_2 a - \lambda_1 b| & \text{if } a - b \le aN_1 - bN_2 \\ (2) \end{cases} \end{cases}$$

Thus to minimize the maximum window size, the value of MWS from equation (2) should be minimized among all unimodular transformations T that are valid for tiling. In many cases, MWS is minimized when  $|\lambda_2 a - \lambda_1 b|$  is minimized.

# 4.2 Legal Transformation

Let  $\vec{d_i} = (d_{i,1}, d_{i,2})$  (i = 1, ..., m) be a set of dependence distance vectors. With uniformly generated references, all the dependences in a nested loop are distance vectors. Given any two uniformly generated references  $\lambda_1 i + \lambda_2 j + c_1$  and  $\lambda_1 i + \lambda_2 j + c_2$ , to test for a dependence from iteration  $(i_1, i_2)$  to iteration  $(j_1, j_2)$ , we check for integer solutions within the loop range to the equation:

$$\lambda_1 i_1 + \lambda_2 i_2 + c_1 = \lambda_1 j_1 + \lambda_2 j_2 + c_2$$

i.e.,

$$\lambda_1(j_1 - i_1) + \lambda_2(j_2 - i_2) = c_1 - c_2$$

We can write  $x_1 = j_1 - i_1$  and  $x_2 = j_2 - i_2$  where  $(x_1, x_2)$  is a distance vector. Since  $\lambda_1, \lambda_2, c_1, c_2$  are constants, every solution gives a distance vector. The smallest lexicographically positive solution is the dependence vector of interest. In order for the transformation T to render the loop nest tileable, the following conditions must hold:

$$ad_{i,1} + bd_{i,2} \ge 0$$
  $i = 1, \cdots, m$   
 $cd_{i,1} + dd_{i,2} \ge 0$   $i = 1, \cdots, m$ 

We illustrate the use of technique through Example 2. Consider the loop nest:

for 
$$i = 1$$
 to 25 do

for 
$$j = 1$$
 to 10 do

X[2i + 5j + 1] = X[2i + 5j + 5]

The distance vectors for this loop are: (3, -2), (2, 0), (5, -2). The problem here is to find a unimodular transformation  $T: T = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$  such that the loop is tileable (which allows bringing chunks of data which can fully operated upon before discarding), *i.e.*, represented by the following constraints:

• 
$$3a - 2b \ge 0, 2a \ge 0, 5a - 2b \ge 0, 3c - 2d \ge 0, 2c \ge 0, 5c - 2d \ge 0.$$

and the maximum window size (MWS)

$$MWS = \begin{cases} \left( \left| \frac{24}{b} \right| + 1 \right) |5a - 2b| & \text{if } a - b \ge 25a - 10b \\ \left( \left| \frac{9}{a} \right| + 1 \right) |5a - 2b| & \text{if } a - b \le 25a - 10b \end{cases}$$

is minimized. Given the set of inequalities that should be satisfied,

$$3a - 2b \ge 0 \Longrightarrow b \le \frac{3a}{2} \Longrightarrow 9b \le \frac{27a}{2}$$

Since,  $9b \leq \frac{27a}{2}$ , the second condition applies, *i.e.*,  $9b \leq 24a$ . So,

$$MWS = \left(\frac{9}{a} + 1\right)(5a - 2b) = 45 + (5a - 2b) - \frac{18b}{a}$$

needs to be minimized subject to inequalities (2.5-2.10). We use either a branch and bound technique (or general nonlinear programming techniques) to minimize this function; the number of variables is linear in the number of nested loops which is usually very small in practice  $(\leq 4)$  resulting in small solution times. Alternately, if we minimize 5a - 2b subject to constraints (2.5–2.10), we get very good solutions in practice. In the example loop nest, a = 2, b = 3 is an optimal solution, giving an minimum MWS estimate of 22 which is very close to the actual minimum MWS which is 21. In general, the system of inequalities arising legal tiling requirement are combined with either  $\hat{a} - b \leq aN_1 - bN_2$  or with  $\hat{a} - b \geq aN_1 - bN_2$ to form two groups of inequalities; if both the groups have valid solutions, we find the best of these. If only one group has valid solutions, the problem is a lot easier. For the solution a = 2, b = 3, the set of values for c and d which give rise to unimodular T while satisfying tiling legality condition is c = 1, d = 1.

# 4.3 Maximum Window Size for 3-Nested Loops

The window size in 3-nested loops cannot be just derived using the coefficients of the access functions. They are estimated using the dependences, or, in other words the null space vector of the access matrices. The window is a function of the null space vector and the loop limits. The window size is estimated using the largest lexicographic dependence vector since it spans the maximum region in the iteration space. Consider the following loop:

for i = 1 to  $N_1$  do for j = 1 to  $N_2$  do for k = 1 to  $N_3$  do  $\cdots A[\dots, \dots] \cdots$ 

Let  $(d_1, d_2, d_3)$  be the null space vector of the data reference matrix of A. The Maximum Window Size (MWS) is given by

$$\begin{cases} d_1(N_2 - |d_2|)(N_3 - |d_3|) + 1 & \text{if } d_2 \leq 0 \\ d_1(N_2 - |d_2|)(N_3 - |d_3|) + |d_2|(N_3 - |d_3|) + 1 & \text{if } d_2 > 0 \end{cases}$$

Note from the above result the maximum window size can be reduced if  $d_2 = 0$  and further reduced if  $d_1$  is made zero. In other words, we get the best locality if we can find transformation such that the dependences are carried by the inner levels.

code	default	MWS <sub>unopt</sub>	MWS <sub>opt</sub>
2_point	4,096	65 (98.4%)	3 (99.9%)
3_point	1,024	68 (93.3%)	35 (96.5%)
sor	1,024	65 (93.6%)	35 (96.5%)
matmult	768	273 (64.4%)	273 (64.4%)
3step_log	2,064	511 (75.2%)	122 (94.0%)
full search	2,064	252 (87.8%)	60 (97.1%)
rasta_flt	5,152	2,040 (60.4%)	127 (97.5%)
Average Reduction:		81.9%	92.3%

Figure 2: Default and estimated memory requirements.

# Example 10:

for i = 1 to 10 do for j = 1 to 20 do for k = 1 to 30 do  $\cdots A[3i + k, j + k]$ 

 $\cdots A[3i+k, j+k] \cdots$ The access matrix is  $\begin{pmatrix} 3 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix}$ . The reuse vector is (1, 3, -3);

the maximum window size is: MWS = 1(30 - 3)(20 - 3) + 3(30 - 3) = 540. As we can see from the above example, the maximum window size can be reduced if the dependences are carried by inner levels. So, legal transformations which make the inner loop carry dependences can be applied to reduce the window size and thus increase locality. In other words, if the data reference matrix is used as part of the transformation matrix, only the inner most level can be made to carry the dependence. Thus, the maximum window size reduces to one if a transformation matrix T with the first two rows the same as the data reference matrix is used, i.e.,

$$T = \left(\begin{array}{ccc} 3 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 0 \end{array}\right).$$

Our algorithm attempts to find a transformation that increases the level of as many reuse vectors as possible. In the case of Example 10, the reuse vector initially is (1, 3, -3), whose level is 1 since the first non-zero occurs in position 1. After the use of transformation *T*, the reuse vector becomes (0, 0, 1) whose level is 3. We omit the details of the algorithm for lack of space.

# 5 Experimental Results

In order to evaluate the proposed estimation and optimization technique, we tested it using seven codes from image and video processing domains: 2\_point and 3\_point are two-point and three-point stencil codes, respectively; sor is a successive-over-relaxation code; matmult is a matrix-multiply kernel; two different motion estimation codes, 3step\_log and full\_search; and finally, rasta\_flt is a filtering routine from MediaBench [13]

Figure 2 presents our results in columns 2 through 4. The column default gives the normal memory size which is the total number of array elements declared. MWS un opt and MWS opt, on the other hand, give the maximum window sizes (MWS) before and after optimizing the code, respectively. In columns 3 and 4, following each number, within parentheses, we also give the percentage reduction with respect to the corresponding value in the second column. We see from these results that estimating the memory consumption (requirements) of the original (unoptimized) codes indicates a 81.9% saving, and that for the optimized codes brings about an average saving of 92.3%. Note that these savings directly correspond to reduction in the required data memory sizes. We also need to mention that except for rasta\_flt, our estimations were exact. In the rasta\_flt code, our estimation is around 13% higher than the actual memory requirement for both the original and the optimized code.

# 6 Related work

The estimation of the number of references to an array in order to predict cache effectiveness in hierarchical memory machines have been discussed by Ferrante et al. [7] and Gallivan et al. [8]. The image of the iteration space onto the array space to optimize global transfers have been discussed in [8]. A framework for estimating bounds for the number of elements accessed only was given. Ferrante et al. gave exact values for uniformly generated references but did not consider multiple references. Also, for non-uniformly generated references, arbitrary correction factors were given for arriving at lower and upper bounds for the number of distinct references. We present a technique in this paper which gives accurate results for most practical cases and very close bounds where ever necessary. Clauss [3] and Pugh [15] have presented more expensive but exact techniques to count the number of distinct accesses. Researchers from IMEC [1, 4] and Zhao and Malik [20] present techniques that estimate the minimum amount of memory required. They do not address the effect of transformations on the amount of minimum memory required.

Our work on loop transformations for improving data locality bears most similarity to work by Gannon et al. [9, 8] and Eisenbeis et al. [5]. They define the notion of a reference window for each dependence arising from uniformly generated references. Unlike our work, they do not use compound transformations - only interchange and reversal are considered. In addition, the use of a reference window and the resultant need to approximate the combination of these windows results in a loss of precision. Fang and Lu [6] present a method to reduce cache thrashing that requires expensive runtime calculation of which iterations are to be executed by the same processor in order to maximize locality. Ferrante et al. [7] present a formula that estimates the number of distinct references to array elements; their technique does not use dependence information. Kennedy and McKinley [11] discuss a strategy to decide the best ordering of loops for locality; their technique is widely applicable but they do not use compound transformations which can dramatically increase locality. Wolf and Lam [12, 18] develop an algorithm that estimates temporal and spatial reuse of data using localized vector space. Their algorithm combines reuses from multiple references. Their method does not use loop bounds and the estimates used are less precise than the ones presented here. Their method performs an exhaustive search of loop permutations that maximizes locality. Li and Pingali [14] discuss the completion of partial transformations derived from the data access matrix of a loop nest; the rows of the data access matrix are selected subscript functions for various array accesses (excluding constant offsets). While their technique exploits reuse arising from input and output dependences, it does not work well with flow or anti-dependences.

# 7 Summary

Minimizing the amount of memory required is very important for embedded systems. The problem of estimating the minimum amount of memory was recently addressed by Zhao and Malik [20]. In this paper, we presented techniques that (i) quickly and accurately estimates the number of distinct array accesses and the minimum amount of memory in nested loops, and (ii) reduces this number through loop-level transformations. The main abstraction that our technique manipulates is that of data dependence and re-use [19]. Since many compilers that target array-dominated codes maintain some sort of data dependence information, implementing our estimation and optimization strategy involves only a small additional overhead. Our experimental results obtained using a set of seven codes show that the proposed techniques are very accurate, and are capable of reducing the memory consumption significantly through high-level optimizations. Work is in progress to extend our techniques to include the effects of memory layouts of arrays, and to extend the scope of transformations used.

**Acknowledgments** The first author, J. Ramanujam was supported in part by NSF Young Investigator Award CCR-9457768 and by NSF grant CCR-0073800.

#### References

- F. Balasa, F. Catthoor and H. De Man. Background memory area estimation for multi-dimensional signal processing systems. *IEEE Trans.* on VLSI Systems, 3(2):157–172, June 1995.
- [2] F. Catthoor, S. Wuytack, E. D. Greef, F. Balasa, L. Nachtergaele, and A. Vandecappelle. *Custom Memory Management Methodology*. Kluwer Academic Publishers, June 1998.
- [3] P. Clauss. Counting solutions to linear and nonlinear constraints through Ehrhart polynomials: Applications to analyze and transform scientific programs. In Proc. ACM Int. Conf. Supercomp., May 1996.
- [4] H. De Greef, F. Catthoor and H. De Man. Array placement for storage size reduction in embedded multimedia systems. Proc. 11th International Conference on Application-Specific Systems, Architectures and Processors, pages 66–75, July 1997.
- [5] C. Eisenbeis, W. Jalby, D. Windheiser and F. Bodin. A strategy for array management in local memory. *Advances in Languages and Compilers for Parallel Computing*, MIT Press, pp. 130–151, 1991.
- [6] Z. Fang and M. Lu. A solution to cache thrashing problem in RISC based parallel processing systems. *Proc. 1991 International Conf. Parallel Processing*, August 1991.
- [7] J. Ferrante, V. Sarkar and W. Thrash. On Estimating and Enhancing Cache Effectiveness. Proc. 4th Workshop on Languages and Compilers for Parallel Computing, August 1991.
- [8] K. Gallivan, W. Jalby and D. Gannon. On the Problem of Optimizing Data Transfers for Complex Memory Systems. *Proc.* 1988 ACM *International Conference on Supercomputing*, St. Malo, France, pp. 238-253.
- [9] D. Gannon, W. Jalby and K. Gallivan. Strategies for Cache and Local Memory Management by Global Program Transformations. *Journal* of Parallel and Distributed Computing, Vol. 5, No. 5, October 1988, pp. 587-616.
- [10] F. Irigoin and R. Triolet. Supernode Partitioning. Proc. 15th Annual ACM Symp. Principles of Programming Languages, San Diego, CA, Jan. 1988, 319-329.
- [11] K. Kennedy and K. McKinley. Optimizing for parallelism and data locality. Proc. 1992 ACM International Conference on Supercomputing, Washington DC, July 1992.
- [12] M. Lam, E. Rothberg and M. Wolf. The cache performance and optimization of blocked algorithms. *Proc. 4th Intl. Conf. Architectural Support for Programming Languages and Operating Systems*, April 1991.
- [13] C. Lee and M. Potkonjak, and W. H. Mangione-Smith. MediaBench: A tool for evaluating and synthesizing multimedia and communications systems. In Proc. 30th Annual International Symposium on Microarchitecture, pages 330–335, 1997.
- [14] W. Li and K. Pingali. A Singular Loop Transformation Framework Based on Non-singular Matrices. Proc. 5th Workshop on Languages and Compilers for Parallel Computing, August 1992.
- [15] W. Pugh. Counting solutions to Presburger formulas: How and why. Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation, 1994.
- [16] A. Sudarsanam and S. Malik. Simultaneous reference allocation in code generation for dual data memory bank ASIPs. ACM Transactions on Design Automation for Electronic Systems, to appear.
- [17] M. Wolf and M. Lam. A Loop Transformation Theory and an Algorithm to Maximize Parallelism. *IEEE Trans. Parallel and Distributed Syst.*, October 1991.
- [18] M. Wolf and M. Lam. A Data Locality Optimizing Algorithm. Proc. ACM SIGPLAN 91 Conf. Programming Language Design and Implementation, pp. 30–44, June 1991.
- [19] M. Wolfe. High Performance Compilers for Parallel Computing, Addison-Wesley Publishing Company, 1996.
- [20] Y. Zhao and S. Malik. Exact memory size estimation for array computation without loop unrolling. In *Proc. Design Automation Conference*, New Orleans, LA, June 1999.