

Automatic Partitioning of Parallel Loops and Data Arrays for Distributed Shared Memory Multiprocessors*

Anant Agarwal, David Kranz
Laboratory for Computer Science, NE43-624
Massachusetts Institute of Technology
Cambridge, MA 02139
Phone: (617) 253-1448
email: agarwal@mit.edu

Venkat Natarajan
Motorola Cambridge Research Center
Cambridge, MA 02139

Abstract

This paper presents a theoretical framework for automatically partitioning parallel loops to minimize cache coherency traffic on shared-memory multiprocessors. While several previous papers have looked at hyperplane partitioning of iteration spaces to reduce communication traffic, the problem of deriving the *optimal* tiling parameters for minimal communication in loops with general affine index expressions has remained open. Our paper solves this open problem by presenting a method for deriving an optimal hyperparallelepiped tiling of iteration spaces for minimal communication in multiprocessors with caches. We show that the same theoretical framework can also be used to determine optimal tiling parameters for both data and loop partitioning in distributed memory multicomputers. Our framework uses matrices to represent iteration and data space mappings and the notion of uniformly intersecting references to capture temporal locality in array references. We introduce the notion of data footprints to estimate the communication traffic between processors and use linear algebraic methods and lattice theory to compute precisely the size of data footprints. We have implemented this framework in a compiler for Alewife, a distributed shared memory multiprocessor.

1 Introduction

Cache-based multiprocessors are attractive because they seem to allow the programmer to ignore the issues of data partitioning and placement. Because caches dynamically copy data close to where it is needed, repeat references to the same piece of data do not require communication over the network, and hence reduce the need for careful data layout. However, the performance of cache-coherent systems is heavily predicated on the degree of temporal locality in the access patterns of the processor. Loop partitioning for cache-coherent multiprocessors is an effort to increase the percentage of references that hit in the cache.

*A short version of this paper appears in ICPP 1993.

The degree of reuse of data, or conversely, the volume of communication of data, depends both on the algorithm and on the partitioning of work among the processors. (In fact, partitioning of the computation is often considered to be a facet of an algorithm.) For example, it is well known that a matrix multiply computation distributed to the processors by square blocks has a much higher degree of reuse than the matrix multiply distributed by rows or columns.

Loop partitioning can be done by the programmer, by the run time system, or by the compiler. Relegating the partitioning task to the programmer defeats the central purpose of building cache-coherent shared-memory systems. While partitioning can be done at run time (for example, see [1, 2]), it is hard for the run time system to optimize for cache locality because much of the information required to compute communication patterns is either unavailable at run time or expensive to obtain. Thus compile-time partitioning of parallel loops is important.

This paper focuses on the following problem in the context of cache-coherent multiprocessors. Given a program consisting of parallel do loops (of the form shown in Figure 1 in Section 2.1), how do we derive the optimal tile shapes of the iteration-space partitions to minimize the communication traffic between processors. We also indicate how our framework can be used for loop and data partitioning for distributed memory machines, both with and without caches.

1.1 Contributions and Related Work

This paper develops a unified theoretical framework that can be used for loop partitioning in cache-coherent multiprocessors, or for loop and data partitioning in multicomputers with local memory.¹ The central contribution of this paper is a method for deriving an optimal hyperparallelepiped tiling of iteration spaces to minimize communication. The tiling specifies both the shape and size of iteration space tiles. Our framework allows the partitioning of doall loops accessing multiple arrays, where the index expressions in array accesses can be any affine function of the indices.

Our analysis uses the notion of uniformly intersecting references to categorize the references within a loop into classes that will yield cache locality. This notion helps specify precisely the set of references that have substantially overlapping data sets. Overlap produces temporal locality in cache accesses. A similar concept of uniformly generated references has been used in earlier work in the context of *reuse* and iteration space tiling [3, 4].

The notion of data footprints is introduced to capture the combined set of data accesses made by references within each uniformly intersecting class. (The term *footprint* was originally coined by Stone and Thiebaut[5].) Then, an algorithm to compute precisely the total size of the data footprint for a given loop partition is presented. Precisely computing of the size of the set of data elements accessed by a loop tile was itself an important and open problem. While general optimization methods can be applied to minimize the size of the data footprint and derive the corresponding loop partitions, we demonstrate several important special cases where the optimization problem is very simple. The size of data footprints can also be used to guide program transformations to achieve better cache performance in uniprocessors as well.

Although there have been several papers on hyperplane partitioning of iteration spaces, the problem of deriving the optimal hyperparallelepiped tile parameters for general affine index expressions has remained open. For example, Irigoin and Triolet [6] introduce the notion of loop partitioning with multiple hyperplanes which results in hyperparallelepiped tiles. The purpose

¹This paper, however, focuses on loop partitioning, but indicates the modifications necessary for data partitioning.

of tiling in their case is to provide parallelism across tiles, and vector processing and data locality within a tile. They propose a set of basic constraints that should be met by any partitioning and derive the conditions under which the hyperplane partitioning satisfies these constraints.

Although their paper describes useful properties of hyperplane partitioning, it does not address the issue of automatically generating the tile parameters. Careful analysis of the mapping from the iteration space to the data space is very important in automating the partitioning process. Our paper describes an algorithm for automatically computing the partition based on the notion of cumulative footprints, derived from the mapping from iteration space to data space.

Abraham and Hudak [7] considered loop partitioning in multiprocessors with caches. However, they dealt only with index expressions of the form index variable plus a constant. They assumed that the array dimension was equal to the loop nesting and focused on rectangular and hexagonal tiles. Furthermore, the code body was restricted to an update of $A[i, j]$.

Our framework, however, does not place these restrictions on the code body. It is able to handle much more general index expressions, and produce parallelogram partitions if desired. We also show that when Abraham and Hudak’s methods can be applied to a given loop nest, our theoretical framework reproduces their results.

Ramanujam and Sadayappan [8] deal with data partitioning in multicomputers with local memory and use a matrix formulation; their results do not apply to multiprocessors with caches. Their theory produces communication-free hyperplane partitions for loops with affine index expressions when such partitions exist. However, when communication-free partitions do not exist, they can deal only with index expression of the form variable plus a constant offset. They further require the loop dimension to be equal to the loop nesting.

In contrast, our framework is able to discover optimal partitions in cases where communication free partitions are not possible, and we do not restrict the loop nesting to be equal to array dimension. In addition, we show that our framework correctly produces partitions identical to those of Ramanujam and Sadayappan when communication free partitions do exist.

In a recent paper, Anderson and Lam [9] derive communication-free partitions for multicomputers when such partitions exist, and block loops into squares otherwise. Our notion of cumulative footprints allows us to derive optimal partitions even when communication-free partitions do not exist.

Gupta and Banerjee [10] address the problem of automatic data partitioning by analyzing the entire program. Although our paper deals with loop and data partitioning for a single loop only, the following differences in the machine model and the program model lead to problems that are not addressed by Gupta and Banerjee: (1) The data distributions considered by them do not include general hyperparallelepipeds. In order to deal with hyperparallelepipeds, one requires the analysis of communication presented in our paper. (2) Their communication model does not take into account caches. (3) They deal with simple index expressions of the form $c_1 * i + c_2$ and not a general affine function of the loop indices.

Our work complements the work of Wolfe and Lam [3] and Schreiber and Dongarra [11]. Wolfe and Lam derive loop transformations (and tile the iteration space) to improve data locality in multiprocessors with caches. They use matrices to model transformations and use the notion of equivalence classes within the set of uniformly generated references to identify valid loop transformations to improve the degree of temporal and spatial locality within a given loop nest. Schreiber and Dongarra briefly address the problem of deriving optimal hyperparallelepiped iteration space tiles to minimize communication traffic (they refer to it as I/O requirements).

However their work differs from this paper in the following ways: (1) Their machine model does not have a processor cache. (2) The data space corresponding to an array reference and the iteration space are isomorphic. These restrictions make the problem of computing the communication traffic much simpler. Also, one of the main issues addressed by Schreiber and Dongarra is the *atomicity requirement* of the tiles which is related to the dependence vectors and this paper is not concerned with those requirements as it is assumed that the iterations can be executed in parallel.

Ferrante, Sarkar, and Thrash [12] address the problem of estimating the number of cache misses for a nest of loops. This problem is similar to our problem of finding the size of the cumulative footprint, but differs in these ways: (1) We consider a tile in the iteration space and not the entire iteration space; our tiles can be hyperparallelepipeds in general. (2) We partition the references into uniformly intersecting sets, which makes the problem computationally more tractable, since it allows us to deal with only the tile at the origin. (3) Our treatment of coupled subscripts is much simpler, since we look at maximal independent columns, as shown in Section 5.2.

1.2 Overview of the Paper

The rest of this paper is structured as follows. Section 2 states our system model and our program-level assumptions. Section 3 first presents a few examples to illustrate the basic ideas behind loop partitioning; it then discusses the notion of data partitioning, and when it is important. Section 4 develops the theoretical framework for partitioning and presents several additional examples. Section 5 extends the basic framework to handle more general expressions, and Section 6 indicates modifications to the basic framework to handle data partitioning and more general types of systems. The framework for both loop and data partitioning has been implemented in the compiler system for the Alewife multiprocessor. The implementation of our compiler system and a sampling of results is presented in Section 7, and Section 8 concludes the paper.

2 Problem Domain and Assumptions

This paper focuses on the problem of partitioning loops in cache-coherent shared-memory multiprocessors. Partitioning involves deciding which loop iterations will run collectively in a thread of computation. Computing loop partitions involves finding the set of iterations which when run in parallel minimizes the volume of communication generated in the system. This section describes the types of programs currently handled by our framework and the structure of the system assumed by our analysis.

2.1 Program Assumptions

Figure 1 shows the structure of the most general single loop nest that we consider in this paper. The statements in the *loop body* have array references of the form $A[\vec{g}(i_1, i_2, \dots, i_l)]$, where the index function is $\vec{g} : \mathcal{Z}^l \rightarrow \mathcal{Z}^d$, l is the loop nesting and d is the dimension of the array A . We assume that all array references within the loop body are unconditional.

We address the problem of loop and data partitioning for index expressions that are affine

```

Doall (i1=l1:u1, i2=l2:u2, ..., il=l1:u1)
    loop body
EndDoall

```

Figure 1: Structure of a single loop nest

functions of loop indices. In other words, the index function can be expressed as,

$$\vec{g}(\vec{i}) = \vec{i}\mathbf{G} + \vec{a} \quad (1)$$

where \mathbf{G} is a $l \times d$ matrix with integer entries and \vec{a} is an integer constant vector of length d , termed the *offset vector*. Note that \vec{i} , $\vec{g}(\vec{i})$, and \vec{a} are row vectors. We often refer to an array reference by the pair (\mathbf{G}, \vec{a}) . (An example of this function is presented in Section 3). Similar notation has been used in several papers in the past, for example, see [3, 4]. All our vectors and matrices have integer entries unless stated otherwise. We assume that the loop bounds are such that the iteration space is rectangular. However, we note that our methods can still be used to derive reasonable partitions when this condition is not met. Loop indices are assumed to take all integer values between their lower and upper bounds, i.e, the strides are one.

Previous work [7, 8, 13] in this area restricted the arrays in the *loop body* to be of dimension exactly equal to the loop nesting. Abraham and Hudak [7] further restrict the *loop body* to contain only references to a single array; furthermore, all references are restricted to be of the form $A[i_1 + a_1, i_2 + a_2, \dots, i_d + a_d]$ where a_j is an integer constant. Matrix multiplication is a simple example that does not fit these restrictions.

Given P processors, the problem of loop partitioning is to divide the iteration space into P tiles such that the total communication traffic on the network is minimized with the additional constraint that the tiles are of equal size, except at the boundaries of the iteration space. The constraint of equal size partitions is imposed to achieve load balancing. We restrict our discussions to hyperparallelepiped tiles, of which rectangular tiles are a special case.

Like [7, 8, 13], we do not include the effects of synchronization in our framework. Synchronization is handled separately to ensure correct behavior. For example, in the doall loop in Figure 1, one might introduce a barrier synchronization after the loop nest if so desired. We also note that in many cases fine-grain data-level synchronization can be used within a parallel do loop to enforce data dependencies and its cost approximately modeled as slightly more expensive communication than usual [14]. See Appendix B for some details.

2.2 System Model

Our analysis applies to systems whose structure is similar to that shown in Figure 2. The system comprises a set of processors, each with a coherent cache. Cache misses are satisfied by global memory accessed over an interconnection network or a bus. The memory can be implemented as a single monolithic module (as is commonly done in bus-based multiprocessors), or in a distributed fashion as shown in the figure. The memory modules might also be implemented on the processing nodes themselves (data partitioning for locality makes sense only for this case). In all cases, our analysis assumes that the cost of a main memory access is much higher than a cache access, and for loop partitioning, our analysis assumes that the cost of the main memory access is the same no matter where in main memory the data is located.

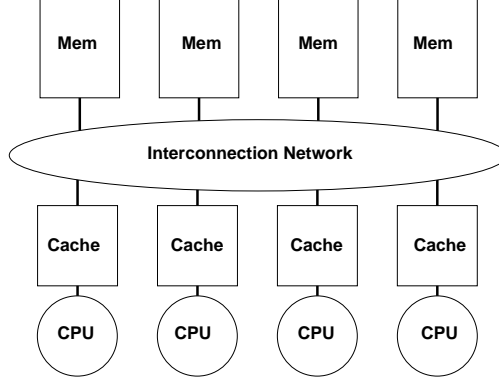


Figure 2: A system with caches and uniform-access main memory.

The goal of loop partitioning is to minimize the total number of main memory accesses. For simplicity, we assume that the caches are large enough to hold all the data required by a loop partition, and that there are no conflicts in the caches. When caches are small, the optimal loop partition does not change, rather, the size of each loop tile executed at any given time on the processor must be adjusted so that the data fits in the cache (if we assume that the cache is effectively flushed between executions of each loop tile). Unless otherwise stated, we assume that cache lines are of unit length. The effect of larger cache lines can be included easily as suggested in [7], and is discussed further in Section 6.2.

3 Loop Partitions and Data Partitions

This section presents examples to introduce and illustrate some of our definitions and to motivate the benefits of optimizing the shapes of loop and data tiles. More precise definitions are presented in the next section.

As mentioned previously, we deal with index expressions that are affine functions of loop indices. In other words, the index function can be expressed as in Equation 1. Consider the following example to illustrate the above expression of index functions.

Example 1 *The reference $A[i_3 + 2, 5, i_2 - 1, 4]$ in a triply nested loop can be expressed by*

$$(i_1, i_2, i_3) \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} + (2, 5, -1, 4)$$

In this example, the second and fourth column of \mathbf{G} are zero indicating that the second and fourth subscripts of the reference are independent of the loop indexes. In such cases, we show in Section 5.2 that we can ignore those columns and treat the referenced array as an array of lower dimension. In future, without loss of generality, we assume that the \mathbf{G} matrix contains no zero columns.

Now, let us introduce the concept of a *loop partition* by examining the following example. Loop partitioning specifies the tiling parameters of the iteration space. Loop partitioning is sometimes termed iteration space partitioning or tiling.

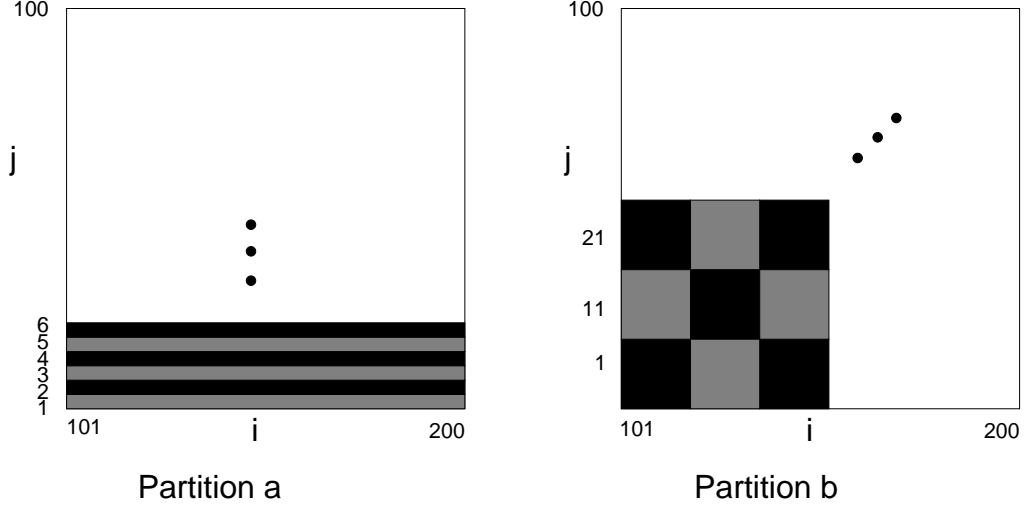


Figure 3: Two simple rectangular loop partitions in the iteration space.

Example 2

```

Doall (i=101:200, j=1:100)
  A[i,j] = B[i+j,i-j-1]+B[i+j+4,i-j+3]
EndDoall

```

Let us assume that we have 100 processors and we want to distribute the work among them. There are 10,000 points in the iteration space and so one can allocate 100 of these to each of the processors to distribute the load uniformly. Figure 3 shows two simple ways of partitioning the iteration space – by rows and by square blocks – into 100 equal tiles.

Minimizing communication volume requires that we minimize the number of data elements accessed by each loop tile. To facilitate this optimization, we introduce the notion of a *data footprint*. Footprints comprise the data elements referenced within a loop tile. In other words, the footprints are regions of the *data space* accessed by a loop tile. In particular, the footprint with respect to a specific reference in a loop tile gives us all the data elements accessed through that reference from within a tile of a loop partition.

Using Figure 4, let us illustrate the footprints corresponding to a reference of the form $B[i+j,i-j]$ for the two loop partitions shown in Figure 3. The footprints in the data space resulting from the loop partition **a** are diagonal stripes and those resulting from partition **b** are square blocks rotated by 45 degrees. Algorithms for deriving the footprints are presented in the next section.

Let us compare the two loop partitions in the context of a system with caches and uniform-access memory (see Figure 2) by computing the number of cache misses. The number of cache misses is equal to the number of distinct elements of **B** accessed by a loop tile, which is equal to the size of a loop tile's footprint on the array **B**. (Section 6.1 deals with minimizing *cache-coherence* traffic). Caches automatically fetch a loop tile's data footprint as the loop tile executes. For each tile in partition **a**, the number of cache misses can be shown to be 104 (see Section 5.1)

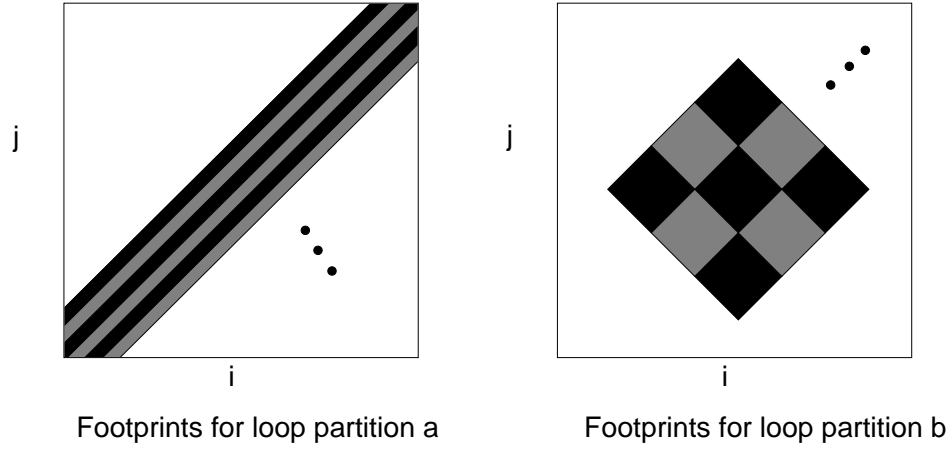


Figure 4: Data footprints in the data space resulting from loop partitions **a** and **b**

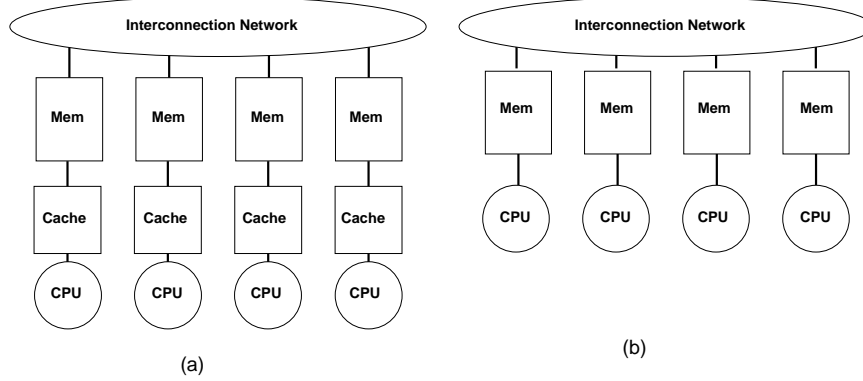


Figure 5: Systems with nonuniform main-memory access time.

whereas the number of cache misses in each tile of partition **b** can be shown to be 140. Thus, because it allows data reuse, loop partition **a** is a better choice if our goal is to minimize the number of cache misses, a fact that is not obvious from the source code.

When is *data partitioning* important? Data partitioning is the problem of partitioning the data arrays into data tiles and assigning each data tile to a local memory module, such that the number of memory references that can be satisfied by the local memory is maximized. Data partitioning is relevant only for nonuniform memory-access (NUMA) systems (for example, see Figure 5).

In systems with nonuniform memory-access times, both loop and data partitioning are required. Our analysis applies to such systems as well. The loop tiles are assigned to the processing nodes and the data tiles to memory modules associated with the processing nodes so that a maximum number of the data references made by the loop tiles are satisfied by the local memory module. Note that in systems with nonuniform memory-access times, but which have caches, data partitioning may still be performed to maximize the number of caches misses that can be satisfied by the memory module local to the processing node.

Referring to Figure 4, the footprint size is minimized by choosing a diagonal striping of the

data space as the data partition, and a corresponding horizontal striping of the iteration space as the loop partition. The additional step of aligning corresponding loop and data tiles on the same node maximizes the number of local memory references.

In fact, the above horizontal partitioning of the loop space and diagonal striping of the data space results in zero communication traffic. Ramanujam and Sadayappan [8] presented algorithms to derive such communication-free partitions when possible. On the other hand, in addition to producing the same partitions when communication-traffic-free partitions exist (see Sections 5.1 and 6.3), our analysis will discover partitions that minimize traffic when such partitions are non-existent as well (see Example 8).

Example 3

```
Doall (i=1:N, j=1:N)
  A[i,j] = B[i,j] + B[i+1,j-2] + B[i-1,j+1]
EndDoall
```

For the loop shown in Example 3, a parallelogram partition results in a lower cost of memory access compared to any rectangular partition since most of the inter iteration communication can be internalized to within a processor for a parallelogram partition (see Section 7.1). Because rectangular partitions often do not minimize communication, we would like to include parallelograms in the formulation of the general loop partitioning problem. In higher dimensions a parallelogram tile generalizes to a hyperparallelepiped; the next section defines it precisely.

4 A Framework for Loop and Data Partitioning

This section first defines precisely the notion of a loop partition and the notion of a footprint of a loop partition with respect to a data reference in the loop. We prove a theorem showing that the number of integer points within a tile is equal to the volume of the tile, which allows us to use volume estimates in deriving the amount of communication. We then present the concept of uniformly intersecting references and a method of computing the cumulative footprint for a set of uniformly intersecting references. We develop a formalism for computing the volume of communication on the interconnection network of a multiprocessor for a given loop partition, and show how loop tiles can be chosen to minimize this traffic. We briefly indicate how the cumulative footprint can be used to derive optimal data partitions for multicomputers with local memory (NUMA machines).

4.1 Loop Tiles in the Iteration Space

Loop partitioning results in a tiling of the iteration space. We consider only hyperparallelepiped partitions in this paper; rectangular partitions are special cases of these. Furthermore, we focus on loop partitioning where the tiles are homogeneous except at the boundaries of the iteration space. Under these conditions of homogeneous tiling, the partitioning is completely defined by specifying the tile at the origin, as indicated in Figure 6. Under homogeneous tiling, the concept

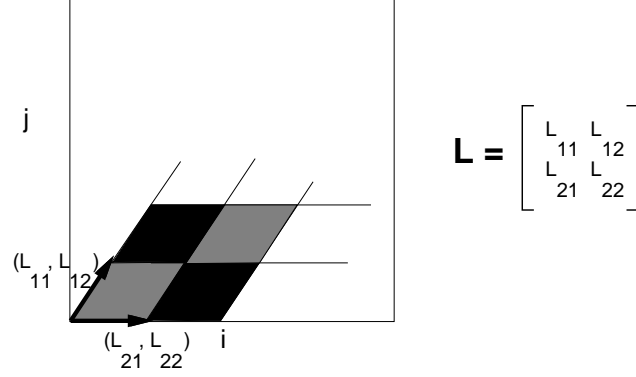


Figure 6: Iteration space partitioning is completely specified by the tile at the origin.

of the tile at the origin is similar to the notion of the clustering basis in [6]. (See Appendix A for a more general representation of hyperparallelepiped loop tiles based on bounding hyperplanes.)

Definition 1 An l dimensional square integer matrix \mathbf{L} defines a semi open hyperparallelepiped tile at the origin of an l dimensional iteration space as follows. The set of iteration points included in the tile is

$$\{\vec{x} \mid \vec{x} = \sum_{i=1}^l \alpha_i \vec{l}_i, 0 \leq \alpha_i < 1\}$$

where \vec{l}_i is the i th row of \mathbf{L} . As depicted in Figure 6, the rows of the matrix \mathbf{L} specify the vertices of the tile at the origin. Often, we also refer to the partition by the \mathbf{L} matrix since each of the other tiles is a translation of the tile at the origin.

Example 4 A rectangular partition can be represented by a diagonal \mathbf{L} matrix. Each row being a separate tile in a three dimensional iteration space $I \times J \times K$ is represented by

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & J & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

Definition 2 A general tile in the iteration space is a translation of the tile at the origin. The translation vector is given by

$$\sum_{i=1}^l \lambda_i \vec{l}_i$$

where λ_i is an integer. A tile is completely specified by $(\lambda_1, \dots, \lambda_l)$. For example $(0, \dots, 0)$ specifies the tile at the origin.

The rest of this paper deals with optimizing the shape of the tile at the origin for minimal communication. Because the amount of communication is related to the number of integer points within a tile, we begin by proving the following theorem relating the volume of a tile to the number of integer points within it. This theorem on lattices allows us to use volumes of hyperparallelepipeds derived using determinants to determine the amount of communication.

Theorem 1 The number of integer points (iteration points) in tile \mathbf{L} is equal to the volume of the tile, which is given by $|\det \mathbf{L}|$.

Proof: We provide a sketch of the proof; a more detailed proof is given in [15].

It is easy to show that the theorem is true for an n -dimensional semi-open rectangle. For a given n -dimensional semi-open hyperparallelepiped, let its volume be V and let P be the number of integer points in it. It can be shown that one can pack R^n of these hyperparallelepipeds into an n -dimensional rectangle of volume V_R and number of integer points P_R , for any positive integer R , such that both $V_R - R^n V$ and $P_R - R^n P$ grow slower than R^n . In other words,

$$V_R = R^n V + v(R), P_R = R^n P + p(R)$$

where $v(R)$ and $p(R)$ grow slower than R^n . Now subtracting the second equation from the first one, and noting that $V_R = P_R$ for the n -dimensional rectangle, we get,

$$V - P = (p(R) - v(R))/R^n.$$

Given that both $v(R)$ and $p(R)$ grow slower than R^n , this can only be true when $V - P = 0$.

Proposition 1 *The number of integer points in any general tile is equal to the number of integer points in the tile at the origin.*

Proof: Straight-forward from the definition of a general tile.

In the following discussion, we ignore the effects of the boundaries of the iteration space in computing the number of integer points in a tile. As our interest is in minimizing the communication for a general tile, we can ignore boundary effects.

4.2 Footprints in the Data Space

For a system with caches and uniform access memory, the problem of loop partitioning is to find an optimal matrix \mathbf{L} that minimizes the number of cache misses. The first step is to derive an expression for the number of cache misses for a given tile \mathbf{L} . Because the number of cache misses is related to the number of unique data elements accessed, we introduce the notion of a *footprint* that defines the data elements accessed by a tile. The footprints are regions of the *data space* accessed by a loop tile.

Definition 3 *The **footprint** of a tile of a loop partition with respect to a reference $A[\vec{g}(\vec{i})]$ is the set of all data elements $A[\vec{g}(\vec{i})]$ of A , for \vec{i} an element of the tile.*

The footprint gives us all the data elements accessed through a particular reference from within a tile of a loop partition. Because we consider homogeneous loop tiles, the number of data elements accessed is the same for each loop tile.

We will compute the number of cache misses for the system with caches and uniform access memory to illustrate the use of footprints. The body of the loop may contain references to several variables and we assume that aliasing has been resolved; two references with distinct names do not refer to the same location. Let A_1, A_2, \dots, A_K be references to array A within the loop body, and let $f(A_i)$ be the *footprint* of the loop tile at the origin with respect to the reference A_i and let

$$f(A_1, A_2, \dots, A_K) = \bigcup_{i=1, \dots, K} f(A_i)$$

be the *cumulative footprint* of the tile at the origin. The number of cache misses with respect to the array A is $|f(A_1, A_2, \dots, A_K)|$. Thus, computing the size of the individual footprints and the size of their union is an important part of the loop partitioning problem.

To facilitate computing the size of the union of the footprints we divide the references into multiple disjoint sets. If two footprints are disjoint or mostly disjoint, then the corresponding references are placed in different sets, and the size of the union is simply the sum of the sizes of the two footprints.

However, references whose footprints overlap substantially are placed in the same set. The notion of *uniformly intersecting references* is introduced to specify precisely the idea of “substantial overlap”. Overlap produces temporal locality in cache accesses, and computing the size of the union of their footprints is more complicated.

The notion of uniformly intersecting references is derived from definitions of intersecting references and uniformly generated references.

Definition 4 *Two references $A[\vec{g}_1(\vec{i})]$ and $A[\vec{g}_2(\vec{i})]$ are said to be intersecting if there are two integer vectors \vec{i}_1, \vec{i}_2 such that $\vec{g}_1(\vec{i}_1) = \vec{g}_2(\vec{i}_2)$. For example, $A[i + c1, j + c2]$ and $A[j + c3, i + c4]$ are intersecting, whereas $A[2i]$ and $A[2i + 1]$ are non-intersecting.*

Definition 5 *Two references $A[\vec{g}_1(\vec{i})]$ and $A[\vec{g}_2(\vec{i})]$ are said to be uniformly generated if*

$$g_1(\vec{i}) = \vec{i}\mathbf{G} + \vec{a}_1 \text{ and } g_2(\vec{i}) = \vec{i}\mathbf{G} + \vec{a}_2$$

where \mathbf{G} is a linear transformation and \vec{a}_1 and \vec{a}_2 are integer constants.

The intersection of footprints of two references that are not uniformly generated is often very small. For non-uniformly generated references, although the footprints corresponding to some of the iteration-space tiles might overlap partially, the footprints of others will have no overlap. Since we are interested in the worst-case communication volume between any pair of footprints, we will assume that the total communication generated by two non-uniformly intersecting references is essentially the sum of the individual footprints.

However, the condition that two references are uniformly generated is not sufficient for two references to be intersecting. As a simple example, $A[2i]$ and $A[2i + 1]$ are uniformly generated, but the footprints of the two references do not intersect. For the purpose of locality optimization through loop partitioning, our definition of reuse of array references will combine the concept of uniformly generated arrays and the notion of intersecting array references. This notion is similar to the equivalence classes within uniformly generated references defined in [3].

Definition 6 *Two array references are uniformly intersecting if they are both intersecting and uniformly generated.*

Example 5 *The following sets of references are uniformly intersecting.*

1. $A[i, j], A[i + 1, j - 3], A[i, j + 4]$.
2. $A[2j, 2, i], A[2j - 5, 2, i], A[2j + 3, 2, i]$.

The following pairs are not uniformly intersecting.

1. $A[i, j], A[2i, j]$.
2. $A[i, j], A[2i, 2j]$.
3. $A[j, 2, i], A[j, 3, i]$.
4. $A[2i], A[2i + 1]$.
5. $A[i + 2, 2i + 4], A[i + 5, 2i + 8]$.
6. $A[i, j], B[i, j]$.

Footprints in the data space for a set of uniformly intersecting references are translations of one another, as shown below. The footprint with respect to the reference (\mathbf{G}, \vec{a}_s) is a translation of the footprint with respect to the reference (\mathbf{G}, \vec{a}_r) , where the translation vector is $\vec{a}_s - \vec{a}_r$.

Proposition 2 *Given a loop tile at the origin \mathbf{L} and references $r = (\mathbf{G}, \vec{a}_r)$ and $s = (\mathbf{G}, \vec{a}_s)$ belonging to a uniformly generated set defined by \mathbf{G} , let $f(r)$ denote the footprint of \mathbf{L} with respect to r , and let $f(s)$ denote the footprint of \mathbf{L} with respect to s . Then $f(s)$ is simply a translation of $f(r)$, where each point of $f(s)$ is a translation of a corresponding point of $f(r)$ by an amount given by the vector $(\vec{a}_s - \vec{a}_r)$. In other words,*

$$f(s) = f(r) \dot{+} (\vec{a}_s - \vec{a}_r).$$

This follows directly from the definition of uniformly intersecting references. Recall that an element \vec{i} of the loop tile is mapped by the reference (\mathbf{G}, \vec{a}_r) to data element $\vec{d}_r = \vec{i}\mathbf{G} + \vec{a}_r$, and by the reference (\mathbf{G}, \vec{a}_s) to data element $\vec{d}_s = \vec{i}\mathbf{G} + \vec{a}_s$. The translation vector, $(\vec{d}_s - \vec{d}_r)$, is clearly independent of \vec{i} .

The volume of cache traffic imposed on the network is related to the size of the cumulative footprint. We describe how to compute the size of the cumulative footprint in the following two sections as outlined below.

- First, we discuss how the size of *the footprint for a single reference* within a loop tile can be computed. In general, the size of the footprint with respect to a given reference is not the same as the number of points in the iteration space tile.
- Second, we describe how the size of *the cumulative footprint for a set of uniformly intersecting references* can be computed. The sizes of the cumulative footprints for each of these sets are then summed to produce the size of the cumulative footprint for the loop tile.

4.3 Size of a Footprint for a Single Reference

This section shows how to compute the size of the footprint (with respect to a given reference and a given loop tile \mathbf{L}) efficiently for certain common cases of \mathbf{G} . The general case of \mathbf{G} is dealt with in Section 5. We begin with a simple example to illustrate our approach.

Example 6

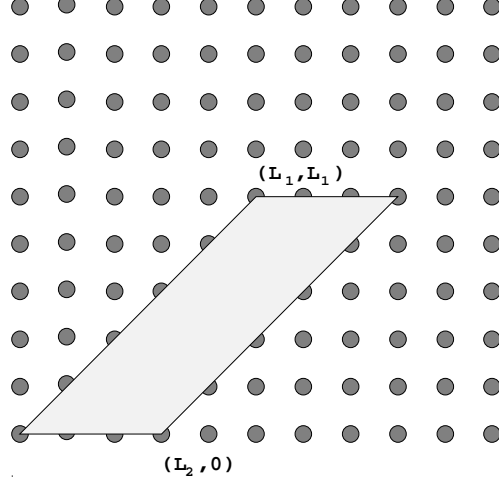


Figure 7: Tile \mathbf{L} at the origin of the iteration space.

```

Doall (i=0:99, j=0:99)
  A[i,j] = B[i+j,j]+B[i+j+1,j+2]
EndDoall

```

The reference matrix \mathbf{G} is

$$\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}.$$

Let us suppose that the loop tile at the origin \mathbf{L} is given by

$$\begin{bmatrix} L_1 & L_1 \\ L_2 & 0 \end{bmatrix}.$$

Figure 7 shows this tile at the origin of the iteration space and the footprint of the tile (at the origin) with respect to the reference $B[i+j, j]$ is shown in Figure 8. The matrix

$$\mathbf{f}(B[i+j, j]) = \mathbf{L}\mathbf{G} = \begin{bmatrix} 2L_1 & L_1 \\ L_2 & 0 \end{bmatrix}$$

describes the footprint. As shown later, the integer points in the semi open parallelogram specified by $\mathbf{L}\mathbf{G}$ is the footprint of the tile and so the size of the footprint is $|\det(\mathbf{L}\mathbf{G})|$. We will use \mathbf{D} to denote the product $\mathbf{L}\mathbf{G}$ as it appears often in our discussion.

The rest of this subsection focuses on deriving the set of conditions under which the footprint size is given by $|\det(\mathbf{D})|$. Briefly, we show that \mathbf{G} being unimodular is a sufficient (but not necessary) condition. The next section derives the size of the cumulative footprint for multiple uniformly intersecting references.

In general, is the footprint exactly the integer points in $\mathbf{D} = \mathbf{L}\mathbf{G}$? If not, how do we compute the footprint? The first question can be expanded into the following two questions.

- Is there a point in the footprint that lies outside the hyperparallelepiped \mathbf{D} ? It follows easily from linear algebra that it is not the case.

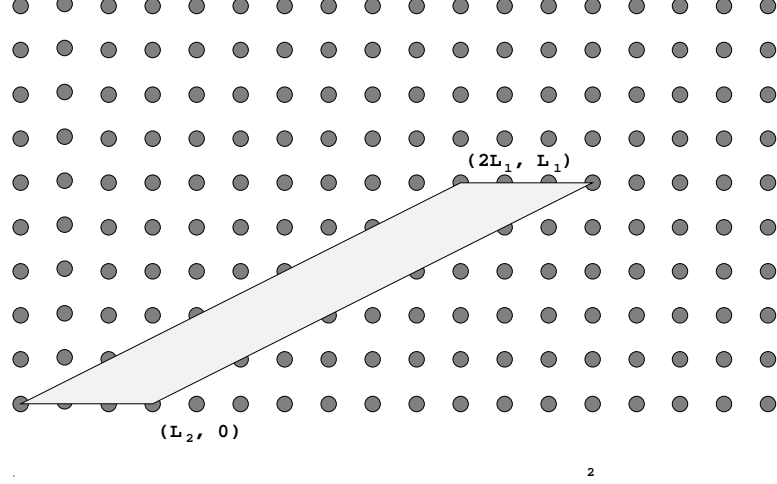


Figure 8: Footprint of \mathbf{L} wrt $B[i+j, j]$ in the data space.

- Is every integer point in \mathbf{D} an element of the footprint? It is easy to show this is not true and a simple example corresponds to the reference $A[2i]$.

We first study the simple case when the hyperparallelepiped \mathbf{D} completely defines the footprint. A precise definition of the set $S(\mathbf{D})$ of points defined by the matrix \mathbf{D} is as follows.

Definition 7 Given a matrix \mathbf{D} whose rows are the vectors \vec{d}_i , $1 \leq i \leq m$, $S(\mathbf{D})$ is defined as the set

$$\{\vec{x} \mid \vec{x} = a_1 \vec{d}_1 + a_2 \vec{d}_2 + \dots + a_m \vec{d}_m, 0 \leq a_i < 1\}.$$

$S(\mathbf{D})$ defines all the points in the semi open hyperparallelepiped defined by \mathbf{D} .

So for the case where \mathbf{D} completely defines the footprint, the footprint is exactly the integer points in $S(\mathbf{D})$. One of the cases where \mathbf{D} completely defines the footprint, is when \mathbf{G} is unimodular as shown below.

Lemma 1 The mapping $\mathbb{R}^l \mapsto \mathbb{R}^d$ as defined by \mathbf{G} is one to one if and only if the rows of \mathbf{G} are independent. Further, the mapping of the iteration space to the data space ($\mathbb{Z}^l \mapsto \mathbb{Z}^d$) as defined by \mathbf{G} is one to one if and only if the rows of \mathbf{G} are independent.

Proof: $\vec{i}_1 \mathbf{G} = \vec{i}_2 \mathbf{G}$ implies $\vec{i}_1 = \vec{i}_2$ if and only if the only solution to $\vec{x} \mathbf{G} = \vec{0}$ is $\vec{0}$. The latter implies that the nullspace of \mathbf{G}^T is of dimension 0. From a fundamental theorem of Linear Algebra [16], this means that the rows of \mathbf{G} are linearly independent. It is to be noted that when the rows of \mathbf{G} are not independent there exists a nontrivial integer solution to $\vec{x} \mathbf{G} = \vec{0}$, given that the entries in \mathbf{G} are integers. This proves the second statement of the lemma.

Lemma 2 The mapping of the iteration space to the data space as defined by \mathbf{G} is onto if and only if the columns of \mathbf{G} are independent and the g.c.d. of the subdeterminants of order equal to the number of columns is 1.

Proof: Follows from the Hermite normal form theorem as shown in [17].

Lemma 3 If \mathbf{G} is invertible then $\vec{d} \in \mathbf{LG}$ if and only if $\vec{d}\mathbf{G}^{-1} \in \mathbf{L}$.

Proof: Clearly \mathbf{G} is invertible implies, $\vec{d}\mathbf{G}^{-1} \in \mathbf{L}$ implies $\vec{d} \in \mathbf{LG}$. \mathbf{G} is invertible implies that the rows of \mathbf{G} are independent and hence the mapping defined by \mathbf{G} is one to one from Lemma 1.

Theorem 2 The footprint of the tile defined by \mathbf{L} with respect to the reference \mathbf{G} is identical to the integer points in the semi open hyperparallelepiped $\mathbf{D} = \mathbf{LG}$ if \mathbf{G} is unimodular.

Proof: It is immediate from Lemma 2 that \mathbf{G} is onto when it is unimodular. \mathbf{G} is onto implies that every data point in \mathbf{D} has an inverse in the iteration space. Can the inverse of the data point be outside of \mathbf{L} ? Lemma 3 shows this is not possible since \mathbf{G} is invertible.

We make the following two observations about Theorem 2.

- \mathbf{G} is unimodular is a sufficient condition; but not necessary. An example corresponds to the reference $A[i + j]$. Further discussions on this is contained in Section 5.
- One may wonder why \mathbf{G} being onto is not sufficient for \mathbf{D} to coincide with the footprint. Even when every integer point in \mathbf{D} has an inverse, it is possible that the inverse is outside of \mathbf{L} . For example, the mapping defined by the \mathbf{G} matrix

$$\begin{bmatrix} 4 \\ 5 \end{bmatrix}$$

is onto as shown by Lemma 2. However, we will show that not all points in \mathbf{LG} are in the footprint. Consider,

$$\mathbf{L} = \begin{bmatrix} 100 & 0 \\ 0 & 100 \end{bmatrix}.$$

The data point (1) is in \mathbf{LG} but none of its inverses is in \mathbf{L} . The same is true for the data points (2), (3), (6), (7), and (11). The one to one property of \mathbf{G} guarantees that no point from outside of \mathbf{L} can be mapped to inside of \mathbf{D} . The reason for this is that the one to one property is true even when \mathbf{G} is treated as a function on reals.

Let us now introduce our technique for computing the cumulative footprint when \mathbf{G} is unimodular. Algorithms for computing the size of the individual footprints and the cumulative footprint when \mathbf{G} is not unimodular are discussed in Section 5.

4.4 Size of the Cumulative Footprint

The size of the cumulative footprint F for a loop tile is computed by summing the sizes of the cumulative footprints for each of the sets of uniformly intersecting references. This section presents a method for computing the size of the cumulative footprint for a set of uniformly intersecting references when \mathbf{G} is unimodular, that is, when the conditions stated in Theorem 2 are true. More general cases of \mathbf{G} are discussed in Section 5. We first describe the method when there are exactly two uniformly intersecting references, and then develop the method for multiple references.

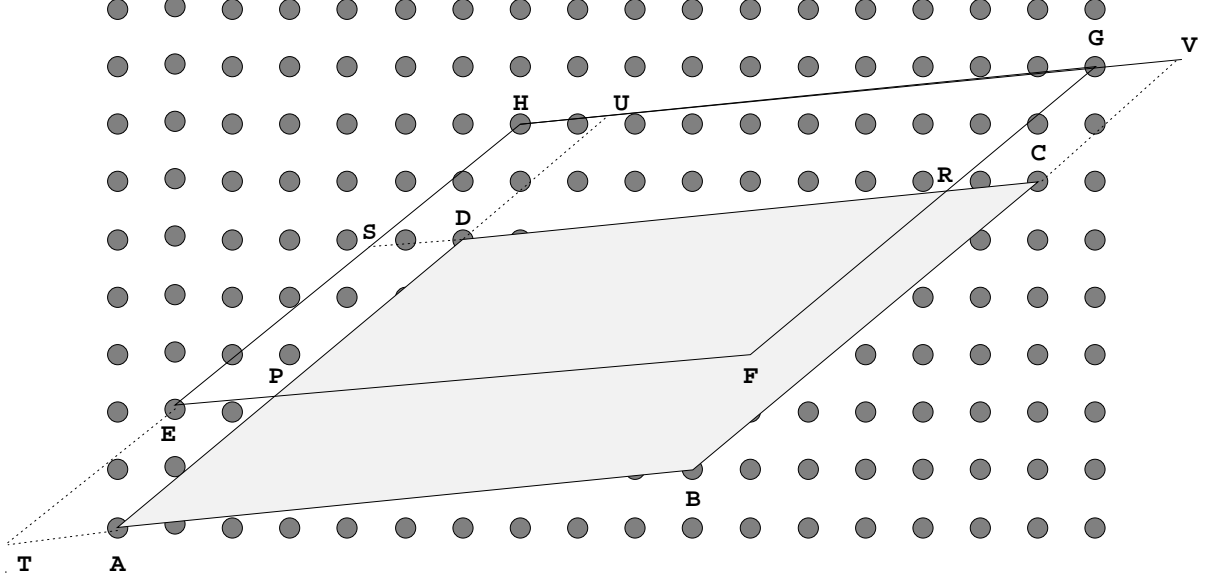


Figure 9: Data footprint wrt $B[i+j, j]$ and $B[i+j+1, j+2]$

Cumulative Footprint for Two References Let us start by illustrating the computation of the cumulative footprint for Example 6. The two references to array B form a uniformly intersecting set and are defined by the following \mathbf{G} matrix.

$$\mathbf{G} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$$

Let us suppose that the loop partition \mathbf{L} is given by

$$\begin{bmatrix} L_{11} & L_{12} \\ L_{21} & L_{22} \end{bmatrix}.$$

Then \mathbf{D} is given by

$$\begin{bmatrix} L_{11} + L_{12} & L_{12} \\ L_{21} + L_{22} & L_{22} \end{bmatrix}.$$

The parallelogram defined by \mathbf{D} in the data space is the parallelogram $ABCD$ shown in Figure 9. $ABCD$ and $EFGH$ shown in Figure 9 are the footprints of the tile \mathbf{L} with respect to the two references ($B[i+j, j]$ and $B[i+j+1, j+2]$ respectively) to array B . In the figure, $\vec{AB} = (L_{11} + L_{12}, L_{12})$, $\vec{AD} = (L_{21} + L_{22}, L_{22})$, and $\vec{AE} = (1, 2)$.

The size of the cumulative footprint is the size of footprint $ABCD$ plus the number of data elements in $EPDS$ plus the number of data elements in $SRGH$. Given that \mathbf{G} is unimodular, the number of data elements is equal to the area $ABCD + SRGH + EPDS = ABCD + ADST + CDUV - SDUH$. Ignoring the area $SDUH$, we can approximate the total area by

$$\left| \det \begin{bmatrix} L_{11} + L_{12} & L_{12} \\ L_{21} + L_{22} & L_{22} \end{bmatrix} \right| + \left| \det \begin{bmatrix} L_{11} + L_{12} & L_{12} \\ 1 & 2 \end{bmatrix} \right| + \left| \det \begin{bmatrix} 1 & 2 \\ L_{21} + L_{22} & L_{22} \end{bmatrix} \right|.$$

Ignoring $SDUH$ is reasonable if we assume that the offset vectors in a uniformly intersecting set of references are small compared to the tile size. We refer to this simplification as the *overlapping*

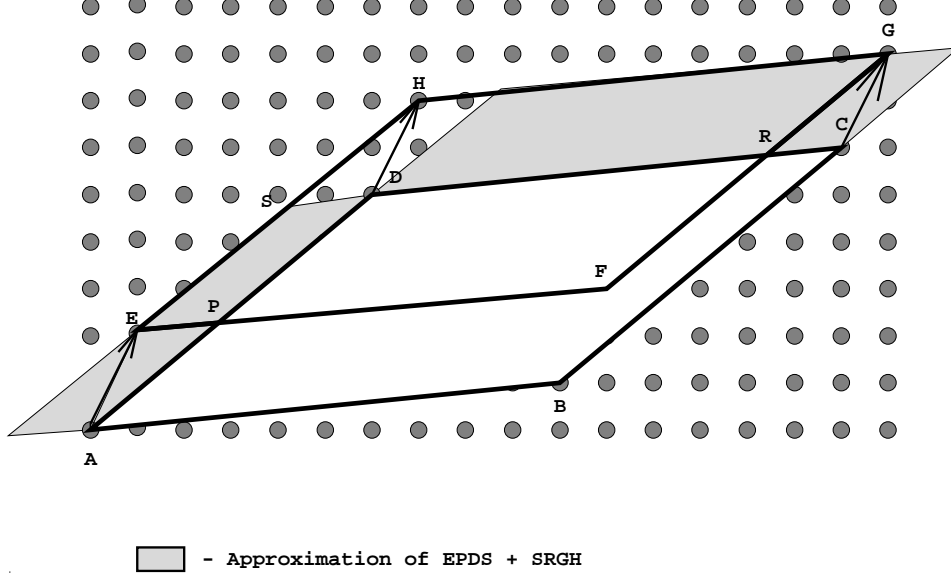


Figure 10: Difference between the cumulative footprint and the footprint.

subtile approximation. This approximation will result in our estimates being higher than the actual values. Although one can easily derive a more exact expression, we use the overlapping subtile approximation to simplify the computation. Figure 16 in Section 7 further demonstrates that the error introduced is insignificant, especially for parallelograms that are near optimal.

The first term in the above equation represents the area of the footprint of a single reference, i.e., $|\det(\mathbf{D})|$. The second and third terms are the determinants of the \mathbf{D} matrix in which one row is replaced by the offset vector $\vec{a} = (1, 2)$. Figure 10 is a pictorial representation of the approximation. The first term is the parallelogram $ABCD$ and the second and third terms are the shaded regions.

The following expression captures the size of the cumulative footprint for the above two references in which one of the offset vectors is $(0, 0)$:

$$|\det \mathbf{D}| + \sum_{k=1}^d |\det \mathbf{D}_{k \rightarrow \vec{a}}|$$

where, $\mathbf{D}_{k \rightarrow \vec{a}}$ is the matrix obtained by replacing the k th row of \mathbf{D} by \vec{a} .

If both the offset vectors are nonzero, because only the relative position of the two footprints determines the area of their nonoverlapping region, we use $\vec{a} = \vec{a}_1 - \vec{a}_0$ in the above equation. The following discussion formalizes this notion and extends it to multiple references.

Cumulative Footprint for Multiple References The basic approach for estimating the cumulative footprint size involves deriving an effective offset vector \hat{a} that captures the combined effects of multiple offset vectors when there are several overlapping footprints resulting from a set of uniformly intersecting references. First, we need a few definitions.

Definition 8 Given a loop tile \mathbf{L} , there are two neighboring loop tiles along the i th row of \mathbf{L} defined by $\{\vec{y} \mid \vec{y} = \vec{x} + \vec{l}_i, \vec{x} \in \text{tile } \mathbf{L}\}$ and $\{\vec{y} \mid \vec{y} = \vec{x} - \vec{l}_i, \vec{x} \in \text{tile } \mathbf{L}\}$, where \vec{l}_i is the i th row

of \mathbf{L} , for $1 \leq i \leq l$. We refer to the former neighbor as the positive neighbor and the latter as the negative neighbor. We also refer to these neighbors as the neighbors of the parallel sides of the tile determined by the rows of \mathbf{L} , excluding the i th row. Figure 11 illustrates the notion of neighboring tiles.

The notion of neighboring tiles can be extended to the data space in like manner as follows.

Definition 9 Given a loop tile \mathbf{L} and a reference (\mathbf{G}, \vec{a}_r) , the neighbors of the data footprint of \mathbf{L} along the k th row of $\mathbf{D} = \mathbf{L}\mathbf{G}$ are $\{\vec{y} | \vec{y} = \vec{x} + \vec{d}_k, \vec{x} \in \mathbf{D} + \vec{a}_r\}$ and $\{\vec{y} | \vec{y} = \vec{x} - \vec{d}_k, \vec{x} \in \mathbf{D} + \vec{a}_r\}$, where \vec{d}_k is the k th row of \mathbf{D} , for $1 \leq k \leq d$.

Definition 10 Given a tile \mathbf{L} , \mathbf{L}' is a subtile wrt the i th row of \mathbf{L} if the rows of \mathbf{L}' are the same as the rows of \mathbf{L} except for the i th row which is α times the i th row of \mathbf{L} , $0 \leq \alpha \leq 1$.

The approximation of the cumulative footprint in Figure 10 can be expressed in terms of subtiles of the tile in the data space. $ABCD$ is a tile in the data space and the two shaded regions in Figure 10 are subtiles of neighboring tiles containing portions of the cumulative footprint. One can view the cumulative footprint as any one of the footprints together with communication from the neighboring footprints. The approximation of the cumulative footprint expresses the communication from the neighboring tiles in terms of subtiles to make the computation simpler.

Definition 11 Let \mathbf{L} be a loop tile at the origin, and let $\vec{g}(\vec{i}) = \vec{i}\mathbf{G} + \vec{a}_r$, $1 \leq r \leq R$ be a set of uniformly intersecting references. For the footprint of \mathbf{L} with respect to any reference (\mathbf{G}, \vec{a}_r) , communication along the positive direction of the k th row of \mathbf{D} is defined as the smallest subtile of the positive neighbor in the k th direction of the footprint which contains the elements of the cumulative footprint within that neighbor. Communication along the negative direction is defined similarly. Communication along the k th row is the sum of these two communications. Each row of \mathbf{D} defines a pair of parallel sides (hyperplanes) of the data footprint determined by the remaining rows of \mathbf{D} . We sometimes refer to the communication along the k th row as the communication across the parallel sides of \mathbf{D} defined by the k th row.

The notion of the communication along the rows of \mathbf{D} facilitates computing the size of the cumulative footprint. Consider the data footprints of a loop tile with respect to a set of uniformly intersecting references shown in Figure 12. The cumulative footprint can be expressed as the union of any one of the footprints and the remaining elements of the cumulative footprint. We take the union because a given data element needs to be fetched only once into a cache.

In Figure 12, the cumulative footprint is the union of the footprint of the loop tile with respect to \vec{a}_4 and the shaded regions corresponding to the remaining elements of the cumulative footprint resulting from the other references. The area of the shaded region can be approximated by the sum of communication along the k th row for $1 \leq k \leq 2$ as shown in Figure 13. The area of the communication along \vec{d}_2 is equal to the area of the parallelogram whose sides are \vec{d}_1 and $\vec{a}_5 - \vec{a}_4$. Among the offset vectors, vector \vec{a}_5 has the maximum component along \vec{d}_2 and vector \vec{a}_4 has the minimum (taking the sign into account) component along \vec{d}_2 . Similarly the area of the communication along \vec{d}_1 is equal to the area of the parallelogram whose sides are \vec{d}_2 and $\vec{a}_4 - \vec{a}_1$ plus the area of the parallelogram whose sides are \vec{d}_2 and $\vec{a}_5 - \vec{a}_4$. This is equal to the area of the parallelogram whose sides are \vec{d}_2 and $\vec{a}_5 - \vec{a}_1$. As before among the offset vectors, vector \vec{a}_5 has the maximum component along \vec{d}_1 and vector \vec{a}_1 has the minimum (taking the

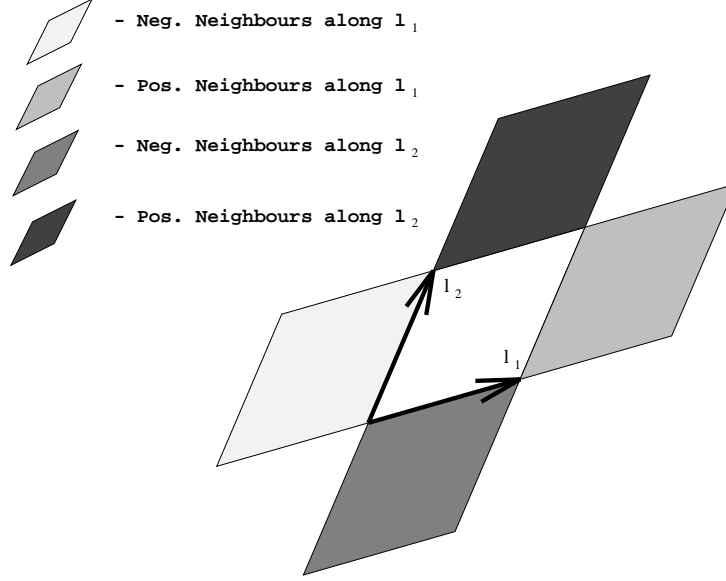


Figure 11: Neighboring Tiles

sign into account) component along \vec{d}_1 . This observation is used in the proof of Theorem 3. It turns out that the effect of offset vector $\vec{a}_5 - \vec{a}_1$ along \vec{d}_2 and $\vec{a}_5 - \vec{a}_4$ along \vec{d}_1 can be captured by a single vector \hat{a} as shown later.

Proposition 3 *Let \mathbf{L} be a loop tile at the origin, and let $\vec{g}(\vec{i}) = \vec{i}\mathbf{G} + \vec{a}_r$, be a set of uniformly intersecting references. The volume of communication along the k th row of \mathbf{D} , $1 \leq k \leq d$, is the same for each of the footprints (corresponding to the different offset vectors).*

Communication along the positive and negative directions will be different for different footprints. But the total communication along the k th row, $1 \leq k \leq d$, is the same for each of the data footprints.

We now derive an expression for the cumulative footprint based on our notion of communication across the sides of the data footprint. Our goal is to capture in a single offset vector \hat{a} the communication in a cache-coherent system resulting from all the offset vectors. More specifically, we would like the k^{th} component of \hat{a} to reflect the communication per unit area across the parallel sides defined by the k th row of \mathbf{D} . The effective vector \hat{a} is derived from the *spread* of a set of offset vectors.

Definition 12 *Given a basis \mathbf{D} and a set of offset vectors \vec{a}_r , $1 \leq r \leq R$, $\text{spread}(\vec{a}_1, \dots, \vec{a}_R)$ is a vector of the same dimension as the offset vectors, whose k th component is given by*

$$\max_r(a_{r,k}) - \min_r(a_{r,k}), \forall k \in 1, \dots, d.$$

In other words, the spread of a set of vectors is a vector in which each component is the difference between the maximum and minimum of the corresponding components in each of the vectors.

The spread as defined above does not quite capture the properties that we are looking for in a single offset vector except when \mathbf{D} is rectangular. To derive the footprint component (or subtile) along a row of \mathbf{D} , we need to compute the difference between the maximum and the

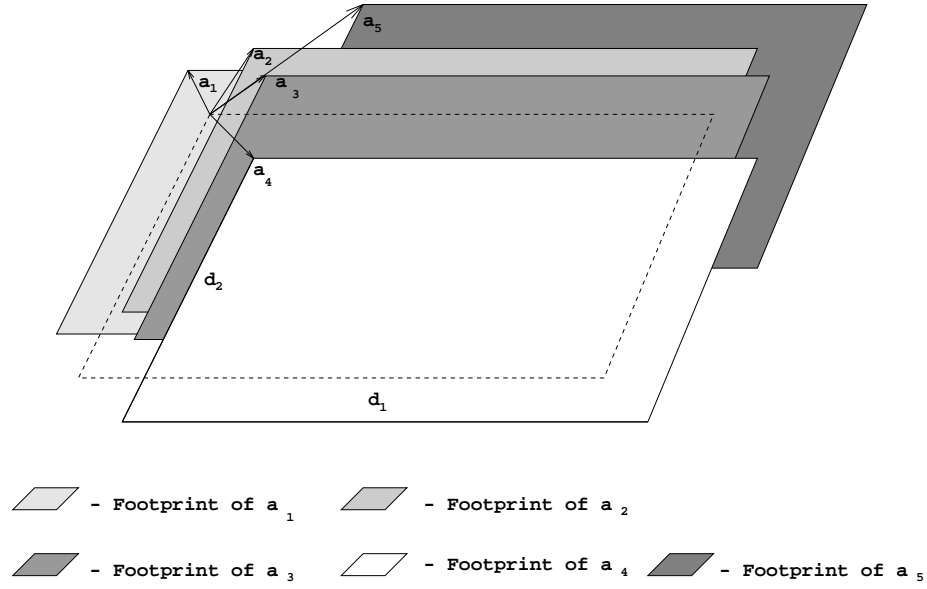


Figure 12: Cumulative Footprint

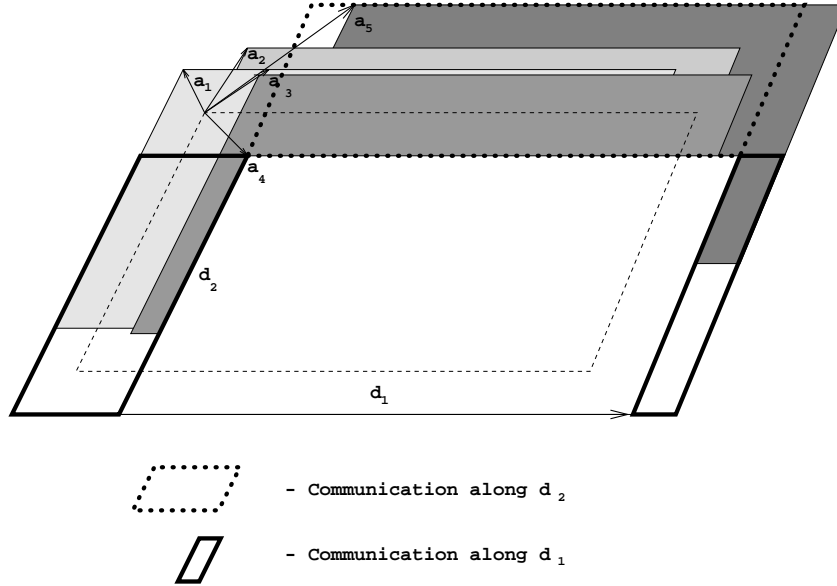


Figure 13: Communication From Neighboring Tiles

minimum components of the offset vectors using \mathbf{D} as a basis. Therefore, we extend the notion of spread to a general basis as follows. Recall that \mathbf{D} is a basis for the data space when \mathbf{G} is unimodular.

In the definition below, \vec{b}_r is the representation of offset vector \vec{a}_r using \mathbf{D} as the basis.

Definition 13 *Given a set of offset vectors \vec{a}_r , $1 \leq r \leq R$, let $\vec{b}_r = \vec{a}_r \mathbf{D}^{-1}$, $\forall r \in 1, \dots, R$ and let \hat{b} be $\text{spread}(\vec{b}_1, \dots, \vec{b}_R)$. Then*

$$\hat{a} = \text{spread}_{\mathbf{D}}(\vec{a}_1, \dots, \vec{a}_R) = \hat{b} \mathbf{D}.$$

Looking at the special case where \mathbf{D} is rectangular helps in understanding the definition.

Proposition 4 *If \mathbf{D} is rectangular then*

$$\hat{a} = \text{spread}(\vec{a}_1, \dots, \vec{a}_R) = \text{spread}_{\mathbf{D}}(\vec{a}_1, \dots, \vec{a}_R)$$

In other words,

$$\hat{a}_k = \max_r(a_{r,k}) - \min_r(a_{r,k}), \forall k \in 1, \dots, d.$$

For example, $\text{spread}_{\mathbf{I}}((1, 0), (2, -1)) = (2 - 1, 0 - 1) = (1, 1)$.

For $\mathbf{D} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$, the spread is given by,

$$\text{spread}_{\mathbf{D}}((1, 0), (2, -1)) = \text{spread}((1, 0)\mathbf{D}^{-1}, (2, -1)\mathbf{D}^{-1})\mathbf{D} = (1, 3)$$

For caches, we use the *max - min* formulation (or the spread) to calculate the amount of communication traffic because the data space points corresponding to the footprints whose offset vectors have values between the *max* and the *min* lie within the cumulative footprint calculated using the spread.²

Lemma 4 *Given a hyperparallelepiped tile \mathbf{L} , and a set of uniformly intersecting references $\vec{g}(\vec{i}) = \vec{i}\mathbf{G} + \vec{a}_r$, where \mathbf{G} is unimodular, the communication along the k th row of $\mathbf{D} = \mathbf{L}\mathbf{G}$ is $\sum_{k=1}^d |\det \mathbf{D}_{k \rightarrow \hat{a}}|$ where $\hat{a} = \text{spread}_{\mathbf{D}}(\vec{a}_1, \dots, \vec{a}_R)$ and $\mathbf{D}_{k \rightarrow \hat{a}}$ is the matrix obtained by replacing the k th row of \mathbf{D} by \hat{a} .*

Proof: *Straight-forward from the definition of spread and the definition of communication along the k th row.*

Theorem 3 *Given a hyperparallelepiped tile \mathbf{L} and a unimodular reference matrix \mathbf{G} , the size of the cumulative footprint with respect to a set of uniformly intersecting references specified by the reference matrix \mathbf{G} and a set of offset vectors $\vec{a}_1, \dots, \vec{a}_R$, is approximately*

$$|\det \mathbf{D}| + \sum_{k=1}^d |\det \mathbf{D}_{k \rightarrow \hat{a}}|$$

where $\hat{a} = \text{spread}_{\mathbf{D}}(\vec{a}_1, \dots, \vec{a}_R)$ and $\mathbf{D}_{k \rightarrow \hat{a}}$ is the matrix obtained by replacing the k th row of \mathbf{D} by \hat{a} .

²For data partitioning, however, the formulation must be modified as discussed in Section 6.3.

Proof: As observed earlier, the size of the cumulative footprint is approximately the size of any of the footprints plus the communication across its sides. Clearly the size of any one of the footprints is given by $|\det \mathbf{D}|$. The rest follows from Lemma 4.

Finally, as stated earlier, the total communication generated by non-uniformly intersecting sets of references is essentially the sum of the communicating generated by the individual cumulative footprints. Example 8 in Section 4.5 discusses an instance of such a computation.

4.5 Minimizing the Size of the Cumulative Footprint

We now focus on the problem of finding the loop partition that minimizes the size of the cumulative footprint. The overall algorithm is summarized in Table 1.

Given:	\mathbf{G} , offset vectors $\vec{a}_1, \dots, \vec{a}_R$
Goal:	Find \mathbf{L} to minimize cumulative footprint size
Procedure:	Write $\mathbf{D} = \mathbf{L}\mathbf{G}$ Find $\vec{b}_1, \dots, \vec{b}_R = \vec{a}_1 D^{-1}, \dots, \vec{a}_R D^{-1}$ Find $\hat{b} = \text{spread}(\vec{b}_1, \dots, \vec{b}_R)$ Then, write $\hat{a} = \hat{b}\mathbf{D}$ Communication $C = \det \mathbf{D} + \sum_{k=1}^d \det \mathbf{D}_{k \rightarrow \hat{a}} $ Finally, find the parameters of \mathbf{L} that minimize C

Table 1: An algorithm for minimizing cumulative footprint size for a single set of uniformly intersecting references. For multiple uniformly intersecting sets, add the communication component due to each set and then determine \mathbf{L} that minimizes the sum.

Let us illustrate this procedure through the following two examples.

Example 7

```

Doall (i=1:N, j=1:N, k=1:N)
  A[i,j,k] = B[i-1,j,k+1] + B[i,j+1,k] + B[i+1,j-2,k-3]
EndDoall

```

Here we have two uniformly intersecting sets of references: one for A and one for B . Let us look at the class corresponding to B since it is more instructive. Because A has only one reference, whose \mathbf{G} is unimodular, its footprint size is independent of the loop partition, given a fixed total size of the loop tile, and therefore need not figure in the optimization process. The G matrix corresponding to the references to B is,

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

The \hat{a} vector is $(2, 3, 4)$. Consider a rectangular partition $\mathbf{L} = \mathbf{A}$ given by

$$\begin{bmatrix} L_i & 0 & 0 \\ 0 & L_j & 0 \\ 0 & 0 & L_k \end{bmatrix}$$

In this example, the \mathbf{D} matrix is the same as the \mathbf{L} matrix. Because \mathbf{D} is rectangular, we can apply Proposition 4 in simplifying the derivation of \hat{a} . The size of the cumulative footprint for B can now be computed according to Theorem 3 as

$$L_i L_j L_k + 2L_j L_k + 3L_i L_k + 4L_i L_j$$

This expression must be minimized keeping $|\det \mathbf{L}|$ (or the product $L_i L_j L_k$) a constant. The product represents the area of the loop tile and must be kept constant to ensure a balanced load. The constant is simply the total area of the iteration space divided by P , the number of processors. For example, if the loop bounds are I , J , and K , then we must minimize $L_i L_j L_k + 2L_j L_k + 3L_i L_k + 4L_i L_j$, subject to the constraint $L_i L_j L_k = IJK/P$.

This optimization problem can be solved using standard methods, for example, using the method of Lagrange multipliers [18]. The size of the cumulative footprint is minimized when L_i , L_j , and L_k are chosen in the proportions 2, 3, and 4, or

$$L_i : L_j : L_k :: 2 : 3 : 4$$

Abraham and Hudak's algorithm [7] gives an identical partition for this example.

We now use an example to show how to minimize the total number of cache misses when there are multiple uniformly intersecting sets of references. The basic idea here is that the references from each set contribute additively to traffic.

Example 8

```

Doall (i=1:N, j=1:N)
  A(i,j) = B(i-2,j) + B(i,j-1) + C(i+j-1,j) + C(i+j+1,j+3)
EndDoall

```

There are three uniformly intersecting classes of references, one for B , one for C , and one for A . Because A has only one reference, its footprint size is independent of the loop partition, given a fixed total size of the loop tile, and therefore need not figure in the optimization process.

For simplicity, let us assume that the tile \mathbf{L} is rectangular and is given by

$$\begin{bmatrix} L_1 & 0 \\ 0 & L_2 \end{bmatrix}.$$

Because \mathbf{G} for the references to array B is the identity matrix, the $\mathbf{D} = \mathbf{L}\mathbf{G}$ matrix corresponding to references to B is same as \mathbf{L} , and the \hat{a} vector is $\text{spread}(-2, 0), (0, -1)) = (2, 1)$. Thus, the size of the corresponding cumulative footprint according to Theorem 3 is

$$\begin{vmatrix} L_1 & 0 \\ 0 & L_2 \end{vmatrix} + \begin{vmatrix} 2 & 1 \\ 0 & L_2 \end{vmatrix} + \begin{vmatrix} L_1 & 0 \\ 2 & 1 \end{vmatrix}.$$

Similarly, \mathbf{D} for array C is

$$\begin{bmatrix} L_1 & 0 \\ L_2 & L_2 \end{bmatrix}.$$

The data footprint \mathbf{D} is not rectangular even though the loop tile is. Using Definition 13, $\hat{a} = \text{spread}_{\mathbf{D}}((-1, 0), (1, 3)) = (4, 3)$, and the size of the cumulative footprint with respect to C is

$$\begin{vmatrix} L_1 & 0 \\ L_2 & L_2 \end{vmatrix} + \begin{vmatrix} 4 & 3 \\ L_2 & L_2 \end{vmatrix} + \begin{vmatrix} L_1 & 0 \\ 4 & 3 \end{vmatrix}.$$

The problem of minimizing the size of the footprint reduces to finding the elements of \mathbf{L} that minimizes the sum of the two expressions above subject to the constraint the area of the loop tile $|\det \mathbf{L}|$ is a constant to ensure a balanced load. For example, if the loop bounds are I, J , then the constraint is $|\det \mathbf{L}| = IJ/P$, where P is the number of processors.

The total size of the cumulative footprint simplifies to $2L_1L_2 + 4L_1 + 3L_2$. The optimal values for L_1 and L_2 can be shown to satisfy the equation $4L_1 = 3L_2$ using the method of Lagrange multipliers.

5 General Case of \mathbf{G}

This section analyzes the size of the footprint and the cumulative footprint for a general \mathbf{G} , that is, when \mathbf{G} is not restricted to be unimodular. The computation of the size of the footprint is by case analysis on the \mathbf{G} matrix.

5.1 \mathbf{G} is Invertible, but not Unimodular

\mathbf{G} is invertible and not unimodular implies that not every integer point in the hyperparallelepiped \mathbf{D} is an image of an iteration point in \mathbf{L} . A unit cube in the iteration space is mapped to a hyperparallelepiped of volume equal to $|\det \mathbf{G}|$. So the size of the data footprint is $|\det \mathbf{D}|/|\det \mathbf{G}| = |\det \mathbf{L}|$. When \mathbf{G} is invertible the size of the data footprint is exactly the size of the loop tile since the mapping is one to one.

Next, the expression for the size of the cumulative footprint is very similar to the one in Theorem 3, except that the data elements accessed are not dense in the data space. That is, the data space is sparse.

Lemma 5 *Given an iteration space \mathcal{I} , a reference matrix \mathbf{G} , and a hyperparallelepiped \mathbf{D}_1 in the data space, if the vertices of $\mathbf{D}_1\mathbf{G}^{-1}$ are in \mathcal{I} then the number of elements in the intersection of \mathbf{D}_1 and the footprint of \mathcal{I} with respect to \mathbf{G} is $|\det \mathbf{D}_1|/|\det \mathbf{G}|$.*

Proof: Clear if one views $\mathbf{D}_1\mathbf{G}^{-1}$ as the loop tile \mathbf{L} .

Theorem 4 *Given a hyperparallelepiped tile \mathbf{L} , and an invertible reference matrix \mathbf{G} , the size of the cumulative footprint with respect to a set of uniformly intersecting references specified by*

the reference matrix \mathbf{G} and a set of offset vectors $\vec{a}_1, \dots, \vec{a}_R$, is approximately

$$\frac{|\det \mathbf{D}| + \sum_{k=1}^d |\det \mathbf{D}_{k \rightarrow \hat{a}}|}{|\det \mathbf{G}|}$$

where $\hat{a} = \text{spread}(\vec{a}_1, \dots, \vec{a}_R, \mathbf{D})$ and $\mathbf{D}_{k \rightarrow \hat{a}}$ is the matrix obtained by replacing the k th row of \mathbf{D} by \hat{a} .

Proof: Using lemma 5 one can construct a proof similar to that of Theorem 3.

Example 2 (repeated below for convenience) possesses a \mathbf{G} that is invertible, but not unimodular.

```
Doall (i=101:200, j=1:100)
  A[i,j] = B[i+j,i-j-1]+B[i+j+4,i-j+3]
EndDoall
```

For this example, the reference matrix \mathbf{G} corresponding to array \mathbf{B} is

$$\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix},$$

and the offset vectors are

$$\vec{a}_0 = (0, -1) \quad \text{and} \quad \vec{a}_1 = (4, 3)$$

Let us find the optimal rectangular partition \mathbf{L} of the form

$$\begin{bmatrix} L_i & 0 \\ 0 & L_j \end{bmatrix}.$$

The footprint matrix \mathbf{D} is given by

$$\begin{bmatrix} L_i & L_i \\ L_j & -L_j \end{bmatrix}.$$

The offset vectors using \mathbf{D} as a basis are

$$\vec{b}_0 = \vec{a}_0 \mathbf{D}^{-1} = (-1/(2L_i), 1/(2L_j)),$$

$$\vec{b}_1 = \vec{a}_1 \mathbf{D}^{-1} = (7/(2L_i), 1/(2L_j)).$$

The vector $\hat{b} = (4/L_i, 0)$ and the vector

$$\hat{a} = \hat{b} \mathbf{D} = (4, 4)$$

The size of the cumulative footprint according to Theorem 4 is

$$\frac{\begin{vmatrix} L_i & L_i \\ L_j & -L_j \end{vmatrix} + \begin{vmatrix} L_i & L_i \\ 4 & 4 \end{vmatrix} + \begin{vmatrix} 4 & 4 \\ L_j & -L_j \end{vmatrix}}{\begin{vmatrix} 1 & 1 \\ 1 & -1 \end{vmatrix}}$$

which is

$$L_i L_j + 4L_j$$

If we constrain $L_i L_j = 100$ for load balance, we get $L_j = 1$ and $L_i = 100$. This partitioning represents horizontal striping of the iteration space.

5.2 Columns of \mathbf{G} are Dependent and the Rows are Independent

We can apply Theorem 4 to compute the size of a footprint when the columns of \mathbf{G} are dependent, as long as the rows are independent. We derive a \mathbf{G}' from \mathbf{G} by choosing a maximal set of independent columns from \mathbf{G} , such that \mathbf{G}' is invertible. We can then apply Theorem 4 to compute the size of the footprint as shown in the following example.

Example 9 *Consider the reference $A[i, 2i, i + j]$ in a doubly nested loop. The columns of the \mathbf{G} matrix*

$$\begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

are not independent. We choose \mathbf{G}' to be

$$\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}.$$

Now \mathbf{D}' completely specifies the footprint. The size of the footprint equals $|\det \mathbf{D}'|$. If we choose \mathbf{G}' to be

$$\begin{bmatrix} 2 & 1 \\ 0 & 1 \end{bmatrix}$$

then the size of the footprint is $|\det \mathbf{D}'|/2$ for the new \mathbf{D}' since $|\det \mathbf{G}'|$ is now 2. But both expressions evaluate to the same value as one would expect.

5.3 The rows of \mathbf{G} are Dependent

The rows of \mathbf{G} are dependent means that the mapping from the iteration space to the data space is many to one. It is hard to derive an expression for the footprint in general when the rows are dependent. However, we can compute the footprint and the cumulative footprint for many special cases that arise in actual programs. In this section we shall look at the common case where the rows are dependent because one or more of the index variables do not appear in the array reference. We shall illustrate our technique with the matrix multiply program shown in Example 10 below. The notation $\text{1}\$C[i, j]$ means that the read-modify-write of $C[i, j]$ is atomic.

Example 10

```

Doall (i, 0, N)
  Doall (j, 0, N)
    Doall (k, 0, N)
      1$C[i, j] = 1$C[i, j] + A[i, k] + B[k, j]
    EndDoall
  EndDoall
EndDoall

```

The references to the matrices \mathbf{A} , \mathbf{B} and \mathbf{C} belong to separate uniformly intersecting references. So the cumulative footprint is the sum of the footprints of each of the references. We will focus on $\mathbf{A}[\mathbf{i}, \mathbf{k}]$ and footprint computation for the other references are similar. The \mathbf{G} matrix for $\mathbf{A}[\mathbf{i}, \mathbf{k}]$ is

$$\begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix}.$$

We cannot apply our earlier results to compute the footprint since \mathbf{G} is a many to one mapping. However, we can find an invertible \mathbf{G}' such that for every loop tile \mathbf{L} , there is a tile \mathbf{L}' such that the number of elements in footprints $\mathbf{L}\mathbf{G}$ and $\mathbf{L}\mathbf{G}'$ are the same. For the current example, \mathbf{G}' is obtained from \mathbf{G} by deleting the row of zeros, resulting in a two dimensional identity matrix. Similarly \mathbf{L}' is obtained from \mathbf{L} by eliminating the corresponding (second) column of \mathbf{L} . Now, it is easy to show that the number of elements in footprints $\mathbf{L}\mathbf{G}$ and $\mathbf{L}\mathbf{G}'$ are the same by establishing a one-to-one correspondence between the two footprints. Let us use this method to compute the size of the footprint corresponding to the reference $\mathbf{A}[\mathbf{i}, \mathbf{k}]$. Let us assume that \mathbf{L} is rectangular to make the computations simpler. Let \mathbf{L} be

$$\begin{bmatrix} L_i & 0 & 0 \\ 0 & L_j & 0 \\ 0 & 0 & L_k \end{bmatrix}.$$

Now \mathbf{L}' is

$$\begin{bmatrix} L_i & 0 \\ 0 & 0 \\ 0 & L_k \end{bmatrix}.$$

So the size of the footprint is $L_i L_k$. Similarly, one can show that the size of the other two footprints are $L_i L_j$ and $L_j L_k$. The cumulative footprint is $L_i L_k + L_i L_j + L_j L_k$ which is minimized when L_i , L_j and L_k are equal.

6 Other System Environments

This section describes how our framework can be used to solve the partitioning problem in a wide range of systems including those with coherent caches, distributed-memory, and non-unit cache line sizes.

6.1 Coherence-Related Cache Misses

Our analysis presented in the previous section was concerned with minimizing the cumulative footprint size. This process of minimizing the cumulative footprint size not only minimizes the number of first-time cache misses, but the number of coherence-related misses as well. For example, consider the **doall** loop embedded within a sequential loop in Example 11.

Example 11

Doseq ($\mathbf{t}=1:T$)

```

Doall (i=1:N,j=1:N)
  A(i,j) = A(i+1,j)
EndDoall
EndDoseq

```

For this example, we have

$$\mathbf{G} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Let us attempt to minimize the cumulative footprint for a loop partition of the form

$$\mathbf{L} = \begin{bmatrix} L_i & 0 \\ 0 & L_j \end{bmatrix}$$

The cumulative footprint size is given by

$$L_i L_j + L_j$$

In a load-balanced partitioning, $|\det \mathbf{L}| = L_i L_j$ is a constant, so the $L_i L_j$ term drops out of the optimization. The optimization process then attempts to minimize L_j , which is proportional to the volume of cache coherence traffic, as depicted in Figure 14.

Let us focus on regions X, Y and Z in Figure 14(c). As explained in Figure 13, the processor working on the loop tile to which these regions belong (say, processor P_O) shares a portion of its cumulative footprint with processors working on neighboring regions in the data space. Specifically, region Z is a subtile of the positive neighbor and region Y is a subtile shared with its negative neighbor. Region X, however, is completely private to P_O .

Let us consider the situation after the first iteration of the outer sequential loop. Accesses of data elements within region X will hit in the cache, and thereby incur zero communication cost. Data elements in region Z, however, potentially cause misses because the processor working on the positive neighbor might have previously written into those elements, resulting in those elements being invalidated from P_O 's cache. Each of these misses by processor P_O suffers a network round trip because of the need to inform the processor working on its positive neighbor to perform a writeback and then to send the data to processor P_O . Furthermore, if the home memory location for the block is elsewhere, the miss requires an additional network roundtrip. Similarly, in region Y, a write by processor P_O potentially incurs two network round trips as well. The two round trips result from the need to invalidate the data block from the cache of the processor working on the negative neighbor, and then to fetch the blocks into P_O 's cache.

In any case, the coherence traffic is proportional to the area of the shared region Z, which is equal to the area of the shared region Y, and is given by L_j .

6.2 Effect of Cache Line Size

The effect of cache line sizes can be incorporated easily into our analysis. Because large cache lines fetch multiple data words at the cost of a single miss, one data space dimension will be favored by the cache. Without loss of generality, let us assume that the j^{th} dimension of

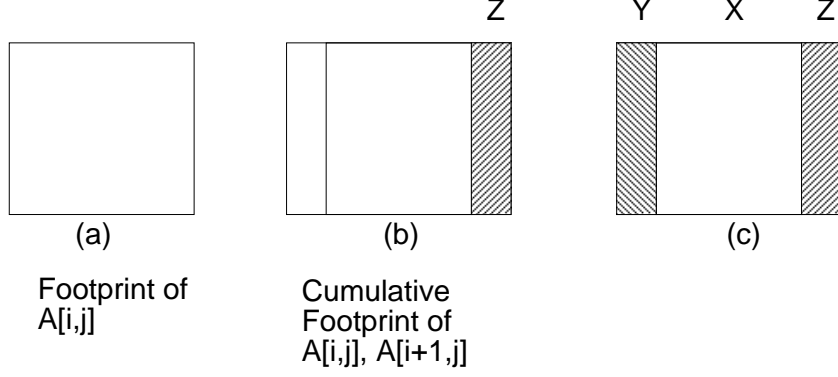


Figure 14: (a) Footprint of reference $A[i, j]$ for a rectangular \mathbf{L} . (b) Cumulative footprint for the references $A[i, j]$ and $A[i + 1, j]$. The hashed region Z represents the increase in footprint size due to the reference $A[i + 1, j]$. (c) The regions X, Y, Z, collectively represent the cumulative footprint for references $A[i, j]$ and $A[i + 1, j]$. Region Z represents the area in the data space shared with the positive neighbor. Region Y represents the area in the data space shared with the negative neighbor.

the data space benefits from larger cache lines. Then, the effect of cache lines of size B can be incorporated into our analysis by replacing each element d_{ij} in the j^{th} column of \mathbf{D} in Theorem 3 by

$$\left\lceil \frac{d_{ij}}{B} \right\rceil$$

to reflect the lower cost of fetching multiple words in the j^{th} dimension of the data space³, and by modifying the definition of intersecting references to the following.

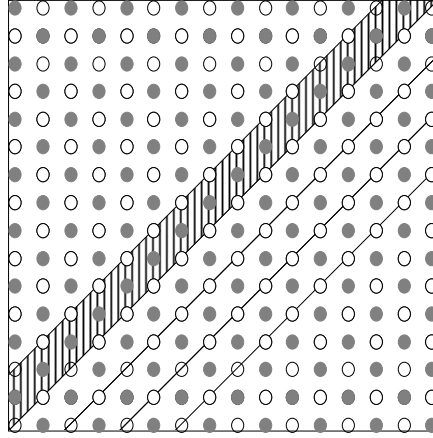
Definition 14 *Two references $A[\vec{g}_1(\vec{i})]$ and $A[\vec{g}_2(\vec{i})]$ are said to be intersecting if there are two integer vectors \vec{i}_1, \vec{i}_2 for which $A[\vec{g}_1(\vec{i}_1)] = A[(d_{11}, d_{12}, \dots)]$ and $A[\vec{g}_2(\vec{i}_2)] = A[(d_{21}, d_{22}, \dots)]$ such that $A[(\dots, d_{1(j-1)}, \lceil \frac{d_{1j}}{B} \rceil, \dots)] = A[(\dots, d_{2(j-1)}, \lceil \frac{d_{2j}}{B} \rceil, \dots)]$, where B is the size of a cache line, and the j^{th} dimension in the data space benefits from larger cache lines.*

6.3 Data Partitioning

In systems in which main memory is distributed with the processing nodes (e.g., see Figure 5), data partitioning is the problem of partitioning the data arrays into data tiles and the nested loops into loop tiles and assigning the loop tiles to the processing nodes and the corresponding data tiles to memory modules associated with the processing nodes so that a maximum number of the data references made by the loop tiles are satisfied by the local memory module. Our formulation facilitates data partitioning straightforwardly. There are two cases to consider: systems with caches and systems without caches.

Systems with Caches The data partitioning strategy in distributed shared-memory systems with caches (Figure 5(a)) proceeds as follows. The optimal loop partition \mathbf{L} is first derived by minimizing the cumulative footprint size as described in the previous sections.

³We note that the estimate of cumulative footprint size will be slightly inaccurate if the footprint is misaligned with the cache block.



Diagonal tiling
of the data space

Figure 15: A communication-free data partition.

Data partitioning requires the additional derivation of the optimal data partition \mathbf{D} for each class of uniformly intersecting references from the optimal loop partition \mathbf{L} . We derive the shapes of the data tiles \mathbf{D} for each \mathbf{G} corresponding to a specific class of uniformly intersecting references. A specific data tile is chosen from the footprints corresponding to each reference in an uniformly intersecting set. In systems with caches, the choice of a specific footprint does not matter, because each data element in the footprint results in a single miss. We then place each loop tile with the data tiles accessed by it on the same processing node.

As an example, let us work out the optimal data partitioning for Example 2. The optimal loop partition for this example was worked out in Section 5.1. The optimal \mathbf{L} was shown to stripe the iteration space horizontally and was given by

$$\begin{bmatrix} 100 & 0 \\ 0 & 1 \end{bmatrix}$$

The corresponding footprint $\mathbf{D} = \mathbf{L}\mathbf{G}$ represents a diagonal striping of the data space and is given by

$$\begin{bmatrix} 100 & 100 \\ 1 & -1 \end{bmatrix}$$

Thus, for this example, if diagonal tiles of data (as depicted in Figure 15) are placed in the memory modules close to the processors with the corresponding iteration tiles, cache misses will be satisfied completely within the node. This data partition thus represents a communication-free data partition.

Interestingly, because \mathbf{G} for this example is not unimodular (its determinant is 2), not all data space points are accessed. In the figure, the shaded points represent the untouched data elements.

Systems without Caches The compiler has two options to optimize communication volume in systems without caches. The compiler can choose to make *local copies* of remote data, or it can fetch remote data each time the data is needed. In the former case, the compiler can use the same partitioning algorithms described in this paper for systems with caches, but it must also solve the data coherence problem for the copied data. This section addresses the latter case.

Although the overall data partitioning strategy remains largely the same as described in the previous section, we must make one change in the footprint size computation to reflect the fact that a given data tile is placed in local memory and data elements from neighboring tiles have to be fetched from remote memory modules each time they are accessed. Because data partitioning for distributed-memory systems without caches (see Figure 5(b)) assumes that data from other memory modules is not dynamically copied locally (as in systems with caches), we replace the *max-min* formulation by the *cumulative spread* a^+ of a set of uniformly intersecting references. That is

$$a^+ = \text{cumulativespread}_{\mathbf{D}}(\vec{a}_1, \dots, \vec{a}_R) = b^+ \mathbf{D},$$

in which the k^{th} element of b^+ is given by,

$$b_k^+ = \sum_r | [b_{r,k} - \text{med}_r(b_{r,k})] |, \forall k \in 1, \dots, d,$$

where $\vec{b}_r = \vec{a}_r \mathbf{D}^{-1}, \forall r \in 1, \dots, R$ and $\text{med}_r(b_{r,k})$ is the median of the offsets in the k^{th} dimension. The rest of our framework for minimizing the footprint size applies to data partitioning if \hat{a} is replaced by a^+ .

The data partitioning strategy proceeds as follows. As in loop partitioning for caches, for a given loop tile \mathbf{L} , we first write an expression for the communication volume by deriving the size of that portion of the cumulative footprint not contained in local memory. This communication volume is given by

$$\sum_{k=1}^d |\det \mathbf{D}_{k \rightarrow a^+}|$$

We then derive the optimal \mathbf{L} to minimize this communication volume. We then derive the optimal data partition \mathbf{D} for each class of uniformly intersecting references from the optimal loop partition \mathbf{L} as described in the previous section on systems with caches. A specific data tile is chosen from the footprints corresponding to each reference in an uniformly intersecting set. In systems without caches, because a single data element might have to be fetched multiple times, the choice of a specific data footprint does matter. A simple heuristic to maximize the number of local accesses is to choose a data tile whose offsets are the medians of all the offsets in each dimension. We can show that using a median tile is optimal for one-dimensional data spaces, and close to optimal for higher dimensions. However, a detailed description is beyond the scope of this paper. We then place each loop tile with the corresponding data tiles accessed by it on the same processor.

7 Implementation and Results

This paper presents cumulative footprint size measurements from an algorithm simulator and execution time measurements from an actual compiler implementation on a multiprocessor.

7.1 Algorithm Simulator Experiments

We have written a simulator of partitioning algorithms that measures the exact cumulative footprint size for any given hyperparallelepiped partition. The simulator also presents analytically computed footprint sizes using the formulation presented in Theorem 3.

We present in Figure 16 algorithm simulator data showing the communication volume for array **B** in Example 3 (repeated below for convenience) resulting from a large number of loop partitions (with tile size 96) representing both parallelograms and rectangles. The abscissa is labeled by the **L** matrix parameters of the various loop partitions, and the parallelogram shape is also depicted above each histogram bar.

```
Doall (i=1:N, j=1:N)
  A[i,j] = B[i,j] + B[i+1,j-2] + B[i-1,j+1]
EndDoall
```

The example demonstrates that the analytical method yields accurate estimates of cumulative footprint sizes. The estimates are higher than the measured values when the partitions are mismatched with the offset vectors due to the overlapping subtile approximation described in Section 4.4. We can also see that the difference between the optimal parallelogram partition and a poor partition is significant. The differences become even greater if bigger offsets are used. This example also shows that rectangular partitions do not always yield the best partition.

7.2 Implementation on Alewife

We have also implemented some of the ideas from our framework in a compiler for the Alewife machine [19] to understand the extent to which good loop partitioning impacts end application performance, and the extent to which our theory predicts the optimal loop partition. The Alewife machine implements a shared global address space with distributed physical memory and coherent caches. The nodes contain slightly modified SPARC processors and are configured in a 2-dimensional mesh network.

Distributed-memory architectures require three types of related analyses to distribute code and data on to the machine:

Loop Partitioning Each processor must be assigned a set of loop iterations that maximizes reuse of data in caches and achieves good load balance.

Data Partitioning and Alignment Arrays must be distributed among the processors such that memory references that miss in the cache go to the local memory rather than across the network to another node. This is accomplished by partitioning arrays with tile shapes suggested by the **D** matrix, and then aligning corresponding loop and data tiles on the same processor.

Placement In an architecture like Alewife the memory access time depends on the distance between the node making the memory request and the node where the requested data resides. The data partitioning and alignment phases make assignments to virtual processors which must be mapped onto the real machine in order to minimize memory reference latency. This is a smaller effect that may become important in very large machines.

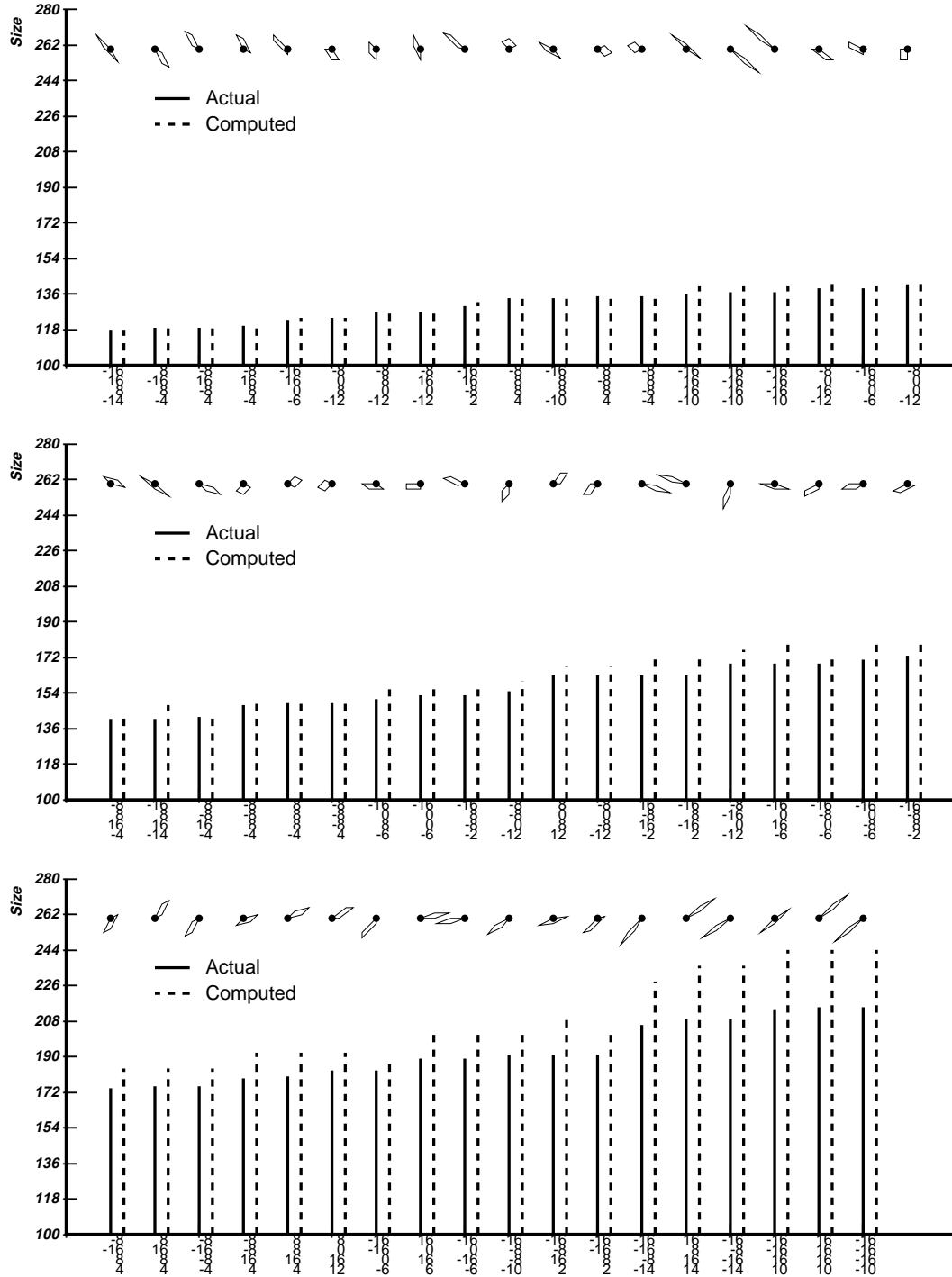


Figure 16: Actual and computed footprint sizes for several loop partitions.

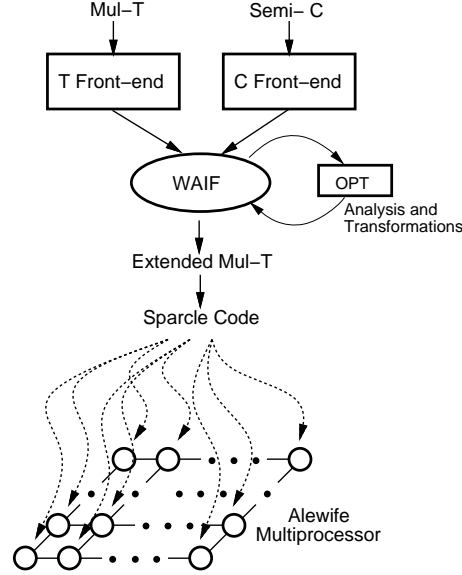


Figure 17: The Alewife Code Generation Process.

We have implemented loop and data partitioning as well as alignment. The results in this paper focus on loop partitioning so we made sure that whatever loop partition was chosen, the optimal data partition for that particular loop partition was used. Otherwise, isolating the effect of cache misses is difficult because changing the loop partition alters both the number of non-local memory references and the number of cache misses.

The structure of our compiler is shown in Figure 17. The input to the compiler is a program where parallelism is specified either by the programmer, or in a previous compilation phase. As in [7], we separate the notion of parallelization from that of implementation. The languages accepted at present are Mul-T, a parallel Lisp language, and Semi-C, a parallel version of C. An initial series of transformations are performed including constant-folding and procedure integration producing a graphical intermediate form called WAIF.

WAIF is a hierarchical graphical representation of a source program. WAIF has two abstraction levels: The program graph (WAIF-PG) and the task and data communications graph (WAIF-CG). WAIF-PG is a customized version of an abstract syntax tree. WAIF-CG summarizes the communication patterns between tasks and data structures that can be derived from a static analysis. Data and loop partitioning are performed as transformations on the WAIF-CG and then code for sequential threads with explicit synchronization is generated. The sequential code-generation process performs standard optimizations such as strength reduction and loop-invariant code motion, producing machine code for Alewife's processors.

7.3 Alewife Experiment

The performance gain due to loop and data partitioning depends on the ratio of communication to computation and other overhead. To get an understanding of these numbers for Alewife, we ran several versions of the following parallel loop nest on an Alewife machine simulator.

```

Doall (i=0:255, j=4:251)
  A[i,j] = A[i-1,j] + B[i,j+4] + B[i,j-4]

```

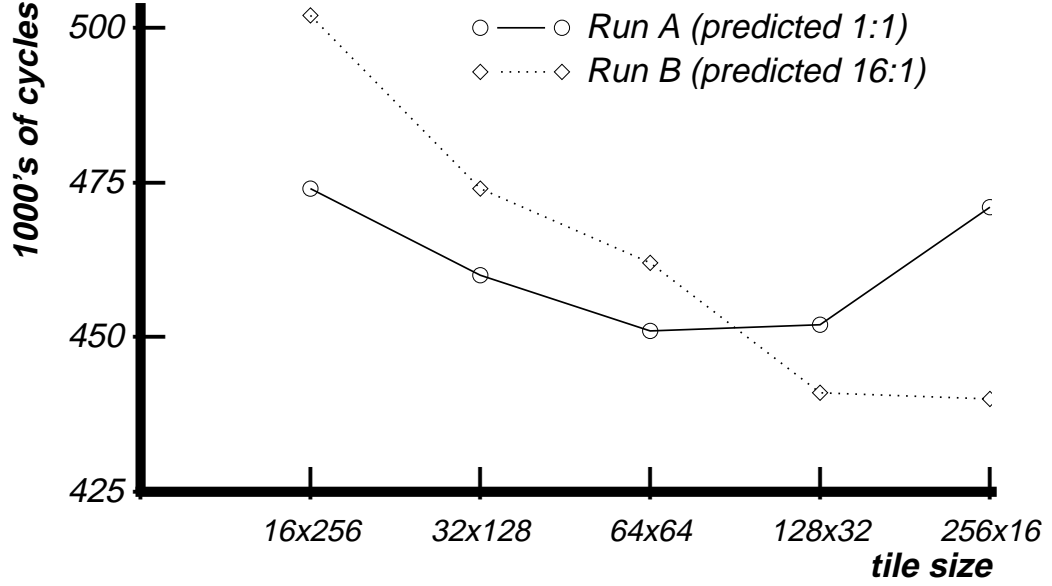


Figure 18: Running times in 1000's of cycles for different aspect ratios on 16 processors.

EndDoall

The \mathbf{G} matrix for the above program fragment is the 2×2 identity matrix, and the offset vectors are $\vec{a}_1 = (0, 0)$, $\vec{a}_2 = (-1, 0)$, $\vec{b}_1 = (0, 4)$, and $\vec{b}_2 = (0, -4)$. We simulated 16 and 64 processors, with each array being 256 elements (words) on a side. The cache line size is four words, and the arrays are stored in row-major order.

Using the algorithms in this paper, and taking the four-word cache line size into account, the compiler chooses a rectangular loop partition and determines that the optimal partition has an aspect ratio of 2:1. The compiler then chooses the closest aspect ratio (1:1) that also achieves load balance for the given problem size and machine size, which results in a tile size of 64x64 iterations. We also generate code using suboptimal partitions with tile sizes ranging from 16x256 to 256x16. This set of executions is labeled run A.

We ran a second version of the program using a different set of offset vectors that give an optimal aspect ratio of 8:1 (run B). This results in a desired tile size between 256x16 and 128x32, with the compiler choosing 256x16, which has the aspect ratio 16:1.

Figure 18 shows the running times for the different tile sizes, and demonstrates that the compiler was able to pick the optimal partitions for both cases. There is some noise in these figures because there can be variation in the cost of accessing the memory that is actually shared due to cache coherence actions, but the minima of the curves are about where the framework predicted. The actual slope of the curves depends on the cost of computing an address and the actual memory latency. One of the reasons that the slopes of these curves are not very steep is that we used ideal data partitions making most references local. The other is that the code generated for index calculations for array references is not very good right now because Alewife does not have virtual memory. This means that arrays must be allocated non-contiguously and accessed indirectly. Much of this overhead can be eliminated with more sophisticated compiler

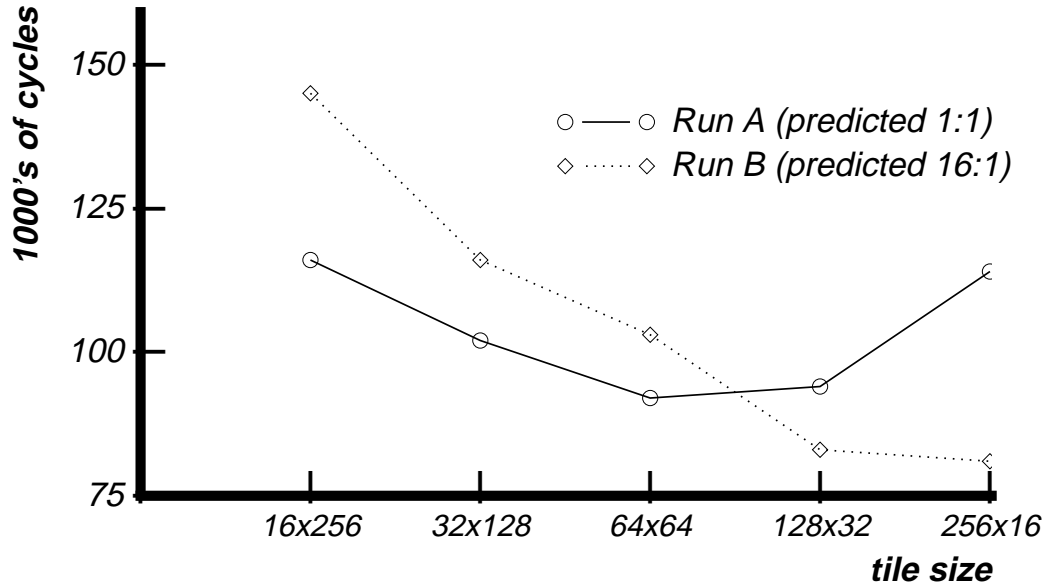


Figure 19: Running times in 1000's of cycles for different aspect ratios with non-communication overhead subtracted out.

analysis of loop invariant expressions.

We can separate these two effects by running the same programs with a no-op replacing the actual loads and stores for the array references. This running time represents the non-communication overhead including index calculation as well as the time to spawn tasks on all processors. Subtracting these times from those in the previous figure gives us the numbers in Figure 19. They represent differences due only to communication and thus represent the greatest possible gain from correct partitioning when the data is partitioned so that almost all accesses are to local memory. In a more realistic program it would likely not be possible to have such an ideal data partition and there would be more non-local references making cache reuse that much more important. In addition, these differences are smaller than they might be on future machines because the local memory latency on Alewife is quite low and will increase as processors get even faster.

Another consideration is that 16 processors is a small machine size. In a larger machine the rectangular partitions can have wider aspect ratios leading to greater differences for non-optimal partitions. We ran the same program on 64 processors, increasing the data size so that each processor would do the same amount of work as in the 16 processor case. The much wider variation can be seen in Figure 20. The 8x512 point for run A is off the top of the chart.

8 Conclusions

The performance of cache-coherent systems is heavily predicated on the degree of temporal locality in the access patterns of the processor. If each block of data is accessed a number of times by a given processor, then caches will be effective in reducing network traffic. Loop partitioning for cache-coherent multiprocessors strives to achieve precisely this goal.

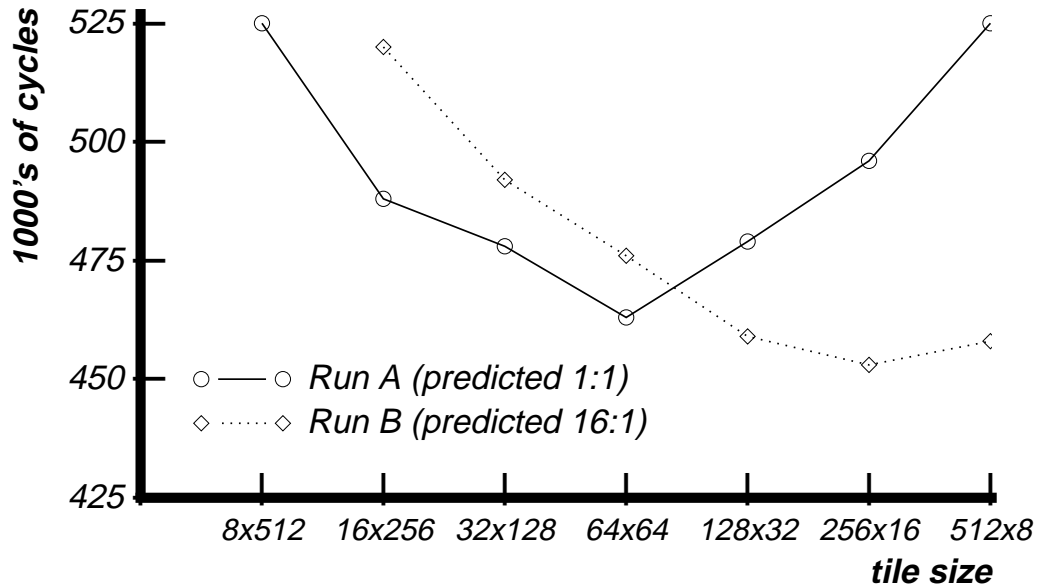


Figure 20: Running times in 1000's of cycles for different aspect ratios on 64 processors.

This paper presented a theoretical framework to derive the parameters of iteration-space partitions of the do loops to minimize the communication traffic in multiprocessors with caches. The framework allows the partitioning of doall loops into optimal hyperparallelepiped tiles where the index expressions in array accesses can be any affine function of the indices. The same framework also yields optimal loop and data partitions for multicomputers with local memory.

Our analysis uses the notion of uniformly intersecting references to categorize the references within a loop into classes that will yield cache locality. The notion of data footprints is introduced to capture the combined set of data accesses made by the references within each uniformly intersecting class. Then, an algorithm to compute precisely the total size of the data footprint for a given loop partition is presented. Once an expression for the total size of the data footprint is obtained, standard optimization techniques can be applied to minimize the size of the data footprint and derive the optimal loop partitions.

Our framework discovers optimal partitions in many more general cases than those handled by previous algorithms. In addition, it correctly reproduces results from loop partitioning algorithms for certain special cases previously proposed by other researchers.

The framework, including both loop and data partitioning for cache-coherent distributed shared memory, has been implemented in the compiler system for the Alewife multiprocessor.

9 Acknowledgments

This research is supported by Motorola Cambridge Research Center and by NSF grant # MIP-9012773. Partial support has also been provided by DARPA contract # N00014-87-K-0825, in part by a NSF Presidential Young Investigator Award. We are grateful to Rajeev Barua for pointing out an error in an earlier formulation of the footprint size and

for extending our compiler implementation to include general affine index expressions and data partitioning. Gino Maa helped define and implement the compiler system and its intermediate form. Andrea Carnevali suggested the simple proof for Theorem 1 on lattices that is sketched in this paper. We acknowledge the contributions of the Alewife group for implementing and supporting the Alewife simulator and runtime system used in obtaining the results.

References

- [1] Constantine D. Polychronopoulos and David J. Kuck. Guided Self-Scheduling: A Practical Scheduling Scheme for Parallel Supercomputers. *IEEE Transactions on Computers*, C-36(12), December 1987.
- [2] E. Mohr, D. Kranz, and R. Halstead. Lazy Task Creation: A Technique for Increasing the Granularity of Parallel Programs. *IEEE Transactions on Parallel and Distributed Systems*, 2(3):264–280, July 1991.
- [3] M. Wolf and M. Lam. A data locality optimizing algorithm. In *Proceedings of the ACM SIGPLAN 91 Conference Programming Language Design and Implementation*, pages 30–44, 1991.
- [4] D. Gannon, W. Jalby, and K. Gallivan. Strategies for cache and local memory management by global program transformation. *Journal of Parallel and Distributed Computing*, 5:587–616, 1988.
- [5] Harold S. Stone and Dominique Thiebaut. Footprints in the Cache. In *Proceedings of ACM SIGMETRICS 1986*, pages 4–8, May 1986.
- [6] F. Irigoin and R. Triolet. Supernode Partitioning. In *15th Symposium on Principles of Programming Languages (POPL XV)*, pages 319–329, January 1988.
- [7] S. G. Abraham and D. E. Hudak. Compile-time partitioning of iterative parallel loops to reduce cache coherency traffic. *IEEE Transactions on Parallel and Distributed Systems*, 2(3):318–328, July 1991.
- [8] J. Ramanujam and P. Sadayappan. Compile-Time Techniques for Data Distribution in Distributed Memory Machines. *IEEE Transactions on Parallel and Distributed Systems*, 2(4):472–482, October 1991.
- [9] Jennifer M. Anderson and Monica S. Lam. Global Optimizations for Parallelism and Locality on Scalable Parallel Machines. In *Proceedings of SIGPLAN '93, Conference on Programming Languages Design and Implementation*, June 1993.
- [10] M. Gupta and P. Banerjee. Demonstration of Automatic Data Partitioning Techniques for Parallelizing Compilers on Multicomputers. *IEEE Transactions on Parallel and Distributed Systems*, 3(2):179–193, March 1992.
- [11] Robert Schreiber and Jack Dongarra. Automatic Blocking of Nested Loops. Technical report, May 1990. RIACS, NASA Ames Research Center, and Oak Ridge National Laboratory.
- [12] J. Ferrante, V. Sarkar, and W. Thrash. *On Estimating and Enhancing Cache Effectiveness*, pages 328–341. Springer-Verlag, August 1991. Lecture Notes in Computer Science: Languages and Compilers for Parallel Computing. Editors U. Banerjee and D. Gelernter and A. Nicolau and D. Padua.
- [13] J. Ramanujam and P. Sadayappan. Tiling multidimensional iteration spaces for nonshared memory machines. In *Proceedings of Supercomputing '91*. IEEE Computer Society Press, 1991.
- [14] G. N. Srinivasa Prasanna, Anant Agarwal, and Bruce R. Musicus. Hierarchical Compilation of Macro Dataflow Graphs for Multiprocessors with Local Memory. To appear in *IEEE Transactions on Parallel and Distributed Systems*. Also available as MIT Laboratory for Computer Science TM-466, June 1992.

- [15] A. Carnevali, V. Natarajan, and A. Agarwal. A Relationship between the Number of Lattice Points within Hyperparallelepipeds and their Volume. In preparation, August 1993, Motorola Cambridge Research Center.
- [16] Gilbert Strang. *Linear algebra and its applications*, volume 3rd edition. Harcourt Brace Jovanovich, San Diego, CA, 1988.
- [17] A. Schrijver. *Theory of Linear and Integer Programming*. John Wiley & Sons, 1990.
- [18] George Arfken. *Mathematical Methods for Physics*. Academic Press, 1985.
- [19] A. Agarwal *et al.* The MIT Alewife Machine: A Large-Scale Distributed-Memory Multiprocessor. In *Proceedings of Workshop on Scalable Shared Memory Multiprocessors*. Kluwer Academic Publishers, 1991. An extended version of this paper has been submitted for publication, and appears as MIT/LCS Memo TM-454, 1991.
- [20] Paul S. Barth, Rishiyur S. Nikhil, and Arvind. M-Structures: Extending a Parallel, Non-strict, Functional Language with State. In *Proceedings of the 5th ACM Conference on Functional Programming Languages and Computer Architecture*, August 1991.
- [21] B.J. Smith. Architecture and Applications of the HEP Multiprocessor Computer System. *Society of Photocopying Instrumentation Engineers*, 298:241–248, 1981.

A A Formulation of Loop Tiles Using Bounding Hyperplanes

A specific hyperparallelepiped loop tile is defined by a set of bounding hyperplanes. Similar formulations have also been used earlier [6].

Definition 15 *Given a l dimensional loop nest \vec{i} , each tile of a hyperparallelepiped loop partition is defined by the hyperplanes given by the rows of the $l \times l$ matrix \mathbf{H} and the column vectors $\vec{\gamma}$ and $\vec{\lambda}$ as follows. The parallel hyperplanes are $\vec{h}_j \vec{i} = \gamma_j$ and $\vec{h}_j \vec{i} = \gamma_j + \lambda_j$, for $1 \leq j \leq l$. An iteration belongs to this tile if it is on or inside the hyperparallelepiped.*

When loop tiles are assumed to be homogeneous except at the boundaries of the iteration space, the partitioning is completely defined by specifying the tile at the origin, namely $(\mathbf{H}, \vec{0}, \vec{\lambda})$, as indicated in Figure 21. For notational convenience, we denote the tile at the origin as \mathbf{L} .

Definition 16 *Given the tile $(\mathbf{H}, \vec{0}, \vec{\lambda})$ at the origin of hyperparallelepiped partition, let $\mathbf{L} = \mathbf{L}(\mathbf{H}) = \Lambda(\mathbf{H}^{-1})^T$, where Λ is a diagonal matrix with $\Lambda_{ii} = \lambda_i$. We refer to the tile by the \mathbf{L} matrix, as \mathbf{L} completely defines the tile at the origin. The rows of \mathbf{L} specify the vertices of the tile at the origin.*

B Synchronization References

Sequential do loops can often be converted to parallel do loops by introducing fine-grain data-level synchronization to enforce data dependencies or mutual exclusion. The cost of synchronization can be approximately modeled as slightly more expensive communication [14]. For example, in the Alewife system the inner loop of matrix multiply can be written using fine-grain synchronization in the form of the loop in Example 12.

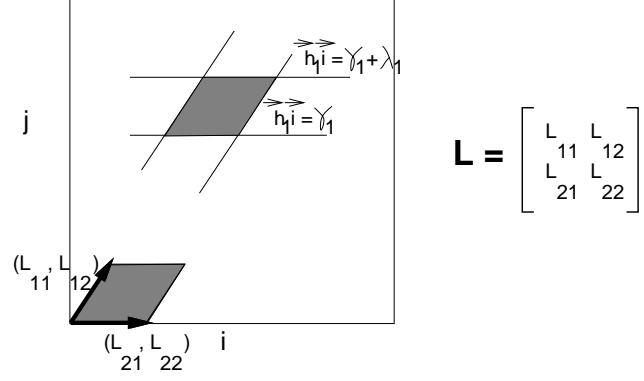


Figure 21: Iteration space partitioning is completely specified by the tile at the origin.

Example 12

```

Doall (i=1:N, j=1:N, k=1:N)
  1$C[i,j] = 1$C[i,j] + A[i,k] + B[k,j]
EndDoall

```

In the code segment in Example 12, the “1\$” preceding the **C** matrix references denote atomic accumulates. Accumulates into the **C** array can happen in any order, just that each accumulate action must be atomic. Such synchronizing reads or writes are both treated as writes by the coherence system. Similar linguistic constructs are also present in Id [20] and in a variant of FORTRAN used on the HEP [21].