

Compiler Support for Array Distribution on NUMA Shared Memory Multiprocessors

TAREK S. ABDELRAHMAN AND THOMAS N. WONG*

tsa@eecg.toronto.edu

*Department of Electrical and Computer Engineering
The University of Toronto
Toronto, Ontario, Canada M5S 1A4*

Editor:

Abstract. Management of program data to improve data locality and reduce false sharing is critical for scaling performance on NUMA shared memory multiprocessors. We use HPF-like data decomposition directives to partition and place arrays in data-parallel applications on Hector, a shared-memory NUMA multiprocessor. We describe a compiler system for automating the partitioning and placement of arrays. The compiler exploits Hector's shared memory architecture to efficiently implement distributed arrays. Experimental results from a prototype implementation demonstrate the effectiveness of these techniques. They also demonstrate the magnitude of the performance improvement attainable when our compiler-based data management schemes are used instead of operating system data management policies; performance improves by up to a factor of 5.

Keywords: data distribution, locality management, cache management, parallelizing compilers, NUMA multiprocessors.

1. Introduction

Shared memory multiprocessors offer a familiar model of programming and better support for operating systems and programming environments. However, bus-based shared memory multiprocessors are severely limited in their scalability. Hence, modern shared memory multiprocessors employ a memory structure in which shared memory is physically distributed, and a scalable point-to-point interconnection network is used to support distributed-shared memory. This structure makes memory access time dependent on the distance between processor and memory. Hence, these multiprocessors are known as Non-Uniform Memory Access (NUMA) multiprocessors. Examples of NUMA multiprocessors include the Stanford Dash [14], the T3D/E systems from Cray [8], the Convex Exemplar SP-1000 [7], and Toronto's Hector multiprocessor [23].

The non-uniform nature of memory access in NUMA multiprocessors makes it necessary to enhance locality of reference in applications. Data must be carefully managed, i.e, partitioned, placed and possibly replicated across the physically distributed memory, to enable a processor to make most of its accesses to the local portion of the shared memory. Such data management tasks are typically delegated to the operating system as part of its page placement policies [12, 13, 20].

* Thomas Wong is presently with IBM Canada. This Work was completed when he was at The University of Toronto.

However, these policies are blind to data access patterns of applications and fail to enhance memory locality over a wide range of applications [1, 12]. Even with full knowledge of data access patterns, operating system policies deal with data in units of pages, which is too coarse of a grain to accommodate the fine-grain sharing that exists in parallel applications. This can result in loss of locality and false sharing. Recent experimental evidence on NUMA systems indicates that significant losses in performance may result from reliance on operating system policies for data management [1, 18].

In this paper we propose and demonstrate that HPF-like directives can be used to partition arrays and place a data partition in the physical memory of the processor using the data the most. More specifically, the paper describes a compiler system for the Hector Multiprocessor which uses data decomposition directives to translate an input data-parallel program, with such directives and parallel loops, into a Single Program Multiple Data (SPMD) program with explicit parallelism and explicitly distributed arrays. The compiler uses array allocation techniques to translate array references in the input program into references to distributed arrays with minimal overhead. Experimental results from a prototype implementation of the compiler demonstrate that improvements in performance can be achieved by using our compiler to manage program data instead of relying on operating system policies.

The remainder of this paper is organized as follows. In section 2, the Hector multiprocessor is briefly described. In section 3, the design of our compiler and the techniques it uses to support array distribution are presented in details. In section 4, experimental results obtained using the prototype implementation of the compiler are presented and discussed. In section 5, related work is reviewed. Finally, in section 6 conclusions and future research directions are presented.

2. The Hector multiprocessor

Hector [23] is a scalable shared memory multiprocessor which consists of a set of processing modules connected by a hierarchical ring interconnection network, as shown in Figure 1. Four processing modules are connected by a bus to form a *station*. Stations are connected by a local ring to form a *cluster*. Clusters are then connected by a central ring to form the system. Each processing module consists of a 20-Mhz Motorola 88000 processor with 16 Kbytes of fast cache and a 4-Mbyte local portion of shared memory. A processing module accesses non-local portions of shared memory over the ring interconnect. A processor accesses its local cache in one processor cycle. In the absence of contention, a processor loads a 16-byte cache line from its local part of the shared memory in 19 processor cycles; from an on-station memory in 29 processor cycles; from an on-ring memory in 37 processor cycles; and from an off-ring memory in 46 processor cycles. There is no hardware mechanism for maintaining cache consistency in Hector; rather consistency is maintained by the operating system.

Hurricane [20] is the operating system of Hector, and it provides application programmers with light weight processes to express and manage parallelism. The

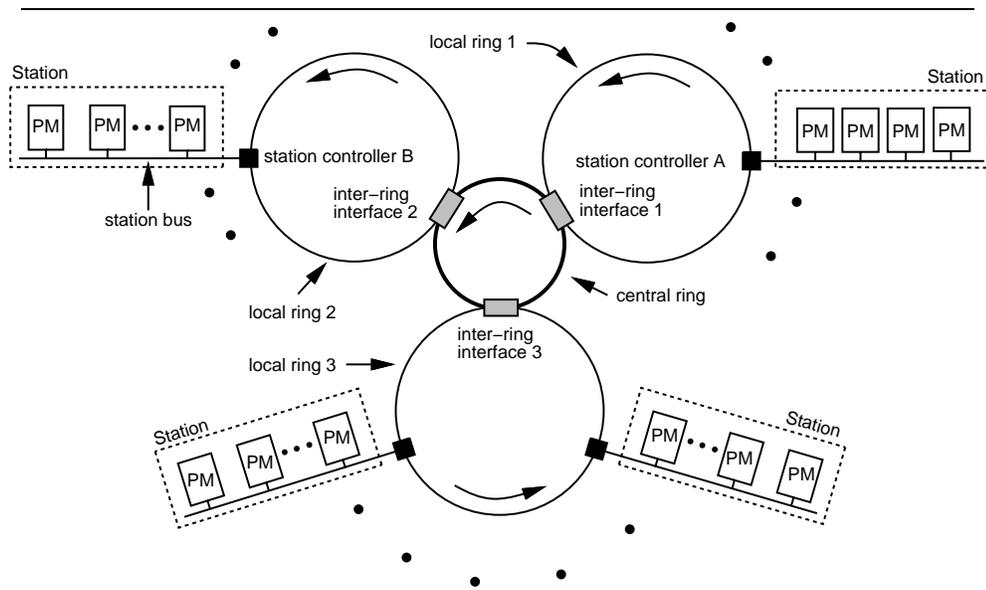


Figure 1. The Hector multiprocessor.

processes run in a single address space and share the same global variables, thus providing programmers with an easy mechanism to share and communicate data.

Hurricane plays a central role in data management on Hector. Since there are no mechanisms for cache consistency in hardware, the operating system maintains the data in the caches in a consistent state. The memory manager of Hurricane implements a single writer multiple reader (SWMR) consistency protocol for each 4-Kbyte page in memory. The manager keeps track of how pages are accessed by the processors. A page is initially cacheable and remains in that state if it is only read, or it is accessed by only one processor. Otherwise, the page becomes uncacheable, and can be accessed only from main memory. Once a page has become uncacheable, it remains as such until the programmer explicitly resets its state back to cacheable. The state-transition diagram of the SWMR protocol is depicted in Figure 2.

The operating system also controls the placement of data in the physically distributed shared memory via two page placement policies. The *First-hit* policy places a page in the physical memory of the first processor that generates a fault to the page. The *round-robin* policy places pages in the physical memories of the system in a round-robin fashion as the processors generate faults to pages. In both cases, once a page has been placed in a memory it remains there until the programmer explicitly resets its state, and then it can be re-placed using one of the above two policies. Pages are not replicated in memory.

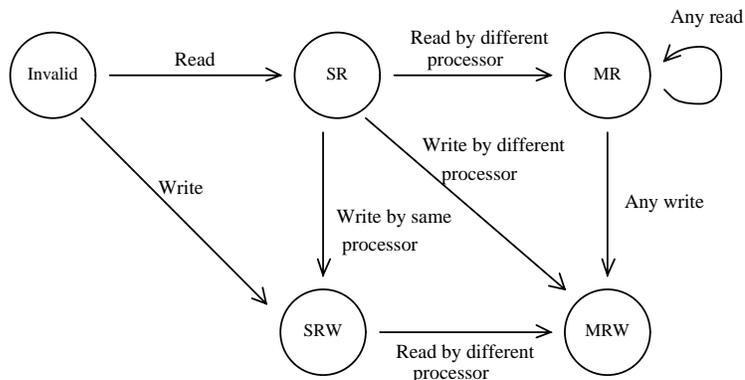


Figure 2. The SWMR coherence protocol.

3. Compiler design

3.1. Overview

The purpose of our compiler is to translate a data parallel program with data decomposition directives and parallel loops into a Single Program Multiple Data (SPMD) program with explicitly distributed arrays and with operating system constructs for the creation and synchronization of parallel tasks. The compiler consists of two main parts: a program analyzer and a program generator. The program analyzer computes the mapping of array partitions to processors, partitions parallel loops and classifies array references made by each processor as either local or non-local to the processor. The program generator produces code to allocate a local region in each physical memory for each array partition that resides in this memory, and produces code to create and synchronize light-weight processes that concurrently execute the iterations of the parallel loops.

3.2. Data decompositions and parallel loops

Data decomposition is achieved by specifying a partitioning scheme for each array in the program and by specifying a processor geometry to which array partitions map. This is done using the `DISTRIBUTE` and `PROCESSOR` directives respectively. The `PROCESSOR` directive declares an n -dimensional Cartesian grid of virtual processors $\mathbf{V}(V_0, V_1, \dots, V_{n-1})$, where V_i is the number of processors in the i^{th} dimension of the grid, and $V_1 \times V_1 \times \dots \times V_{n-1} = P$, the total number of processors. For example, `PROCESSOR P1(32)` declares a linear array P1 of 32 virtual processors. Similarly, `PROCESSOR P2(4,6)` declares a 2-dimensional virtual processor grid P2, with 4 rows and 6 columns.

The `DISTRIBUTE` directive partitions an array by assigning a *partitioning attributes* to each dimension of the array. The `DISTRIBUTE` directive also maps the

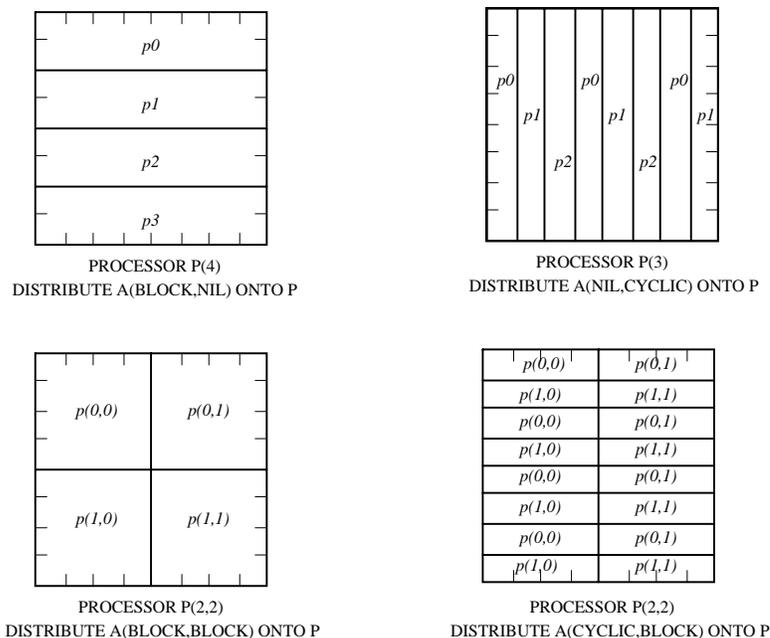


Figure 3. Array decomposition examples.

resulting array partitions onto a processor geometry declared using the PROCESSOR directive. There are four partitioning attributes. The BLOCK attribute divides the corresponding dimension of the array in approximately equal size blocks such that a processor owns a contiguous range of that dimension of the array. The CYCLIC attribute divides the corresponding array dimension by distributing the array elements in this dimension to processors in a round-robin fashion. The BLOCK_CYCLIC attribute first groups array elements in the corresponding dimension in contiguous blocks of a given size, and then assigns the blocks to processors in a round-robin fashion. The block size, called the *block-cyclic factor*, is supplied by the programmer. Finally, the NIL attribute is used to indicate that the corresponding dimension of the array is not distributed. The application of the partitioning attributes to the dimensions of a multi-dimensional array results in useful data decomposition schemes. This is illustrated in Figure 3 using a 2-dimensional array A. The processor geometry on which the array is mapped determines the number of processors assigned to each distributed dimension of the array. For example DISTRIBUTE A(BLOCK,BLOCK) ONTO P, where P is declared as PROCESSOR P(2,4) distributes the array A on 8 processors, assigning 2 processors to the first dimension of A and 4 processors to the second dimension.

Parallel loops are indicated in the input program using FORALL statements. A FORALL statement specifies a loop induction variable, a lower bound, an upper bound and a step size. The use of a FORALL statement implies that: (1) the itera-

tions of the loop may be executed independently in parallel, and (2) all statements in a single iteration are executed sequentially without synchronization.

3.3. The program analyzer

The program analyzer uses the array decomposition directives specified in the input program to perform three main tasks: map array partitions to physical processors, partition parallel loops, i.e., determine the subset of loop iterations that will be executed by each processor, and analyze interprocessor communication to classify array references made by a processor as either local or non-local to the processor. The analysis performed by the program analyzer is largely based on that of Fortran D [9].

The mapping of data partitions to physical processors, is performed in two steps. In the first step, the array partitions are assigned to the virtual processors of the n -dimensional grid to which they are distributed. This is done using the DISTRIBUTE and PROCESSOR directives specified in the input program. In the second step, the virtual processors are assigned to the physical processors. The physical processors are viewed as a linear array, and are numbered $0, 1, \dots, P - 1$. Virtual processor

$(v_0, v_1, \dots, v_{n-1})$ is assigned to the physical processor numbered $\sum_{i=0}^{n-1} v_i \prod_{j=i+1}^{n-1} V_j$.

Thus, for a 2-dimensional grid of virtual processors, this mapping scheme implies a column-major order assignment of virtual processors to physical processors. This mapping scheme does not take into consideration communication patterns in a parallel application, and how to best match them to the network topology on the target multiprocessor to minimize communication cost.

The program analyzer uses the owner-compute rule [9] to partition parallel loops. The owner-compute rule specifies that a processor executes a statement that *computes* an array element if and only if it *owns* the array element, i.e., has the element in its local memory. A processor executes a loop iteration if and only if the execution of this iteration causes the processor to write the array element. Consequently, the subset of loop iterations a processor executes is the subset of loop iterations that write local array elements. The partitioning of a parallel loop is performed by computing the *Local Index Set* for each distributed array on each processor. This is the set of array element indices that are locally owned by a processor. It is calculated using the array decomposition schemes provided by the programmer. The *Global Iteration Set* is computed for each array reference in the left hand side (lhs) of an assignment statement in the parallel loop. This set describes the iteration space spanned by the reference in the parallel loop, and is calculated using the inverse array subscript expressions appearing in the array reference. Finally, the *Local Iteration Set* is computed for the parallel loop on each processor. This set is the subset of loop iterations which are executed by each processor, and is computed by intersecting the Global Iteration Sets with the set that describes the loop iteration space.

```

PROCESSOR P(4,4);
DISTRIBUTE a (BLOCK,BLOCK) ONTO P;
DISTRIBUTE b (BLOCK,BLOCK) ONTO P;
:
:
FORALL (i, 0, 63, 1) /* L1 */
  FORALL (j, 1, 62, 1) /* L2 */
    a[i][j] = ( b[i][j-1] + b[i][j+1] ) / 2

```

(a) A data parallel code segment.

```

Local Index sets of a and b = [32 : 47 : 1][16 : 31 : 1]
Global Iteration Set of a   = [0 : 63 : 1][1 : 62 : 1]
Local Iteration Set of L1   = [32 : 47 : 1]
Local Iteration Set of L2   = [16 : 31 : 1]
Local Access Set of b       = [32 : 47 : 1][15 : 30 : 1] + [32 : 47 : 1][17 : 32 : 1]
Remote Access Set of b      = [32 : 47 : 1][15 : 15 : 1] + [32 : 47 : 1][32 : 32 : 1]
                             + [32 : 47 : 1][15 : 32 : 17]
Remote ownership of b       = { P(2,0) : [32 : 47 : 1][15 : 15 : 1],
                               P(2,2) : [32 : 47 : 1][32 : 32 : 1],
                               others: [NIL] }

```

(b) Index sets for processor P(2,1).

Figure 4. Program analysis for a simple loop nest.

The program analyzer also uses the owner-compute rule to determine which array references are non-local to each processor. Such references must be on the right hand side (rhs) of the statements of the body of the loop since the owner-compute rule makes all references on the lhs local. The *Local Access Set* for each rhs array reference in a loop statement is calculated by determining array indices spanned by the loop's Local Iteration Set. The *Remote Index Set* for each array reference is then computed as the set difference of the Local Access Set and the Local Index Set for the reference. Finally, the ownership of the array elements in the Remote Index Set is determined by intersecting the Remote Index Set with the Local Index Set for the array on each processor.

The analysis performed by the program analyzer is illustrated in Figure 4. It shows a simple example and the various index sets generated by the analyzer on a processor. The index sets are expressed using regular section descriptors, or RSD's [3]. An RSD is denoted by $[l_1 : u_1 : s_1][l_2 : u_2 : s_2] \cdots [l_n : u_n : s_n]$, where l_i , u_i and s_i are respectively the lower bound, upper bound and stride of the i^{th} dimension of the RSD.

3.4. *The program generator*

The program generator produces code to create light-weight processes, or threads, that execute the iterations of the parallel loops concurrently. It replaces loop bounds in the input data parallel program by loop bound variables that are local to each processor. These variables are initialized to different values according to the physical processor number and to the Local Iteration Set of the loop, hence implementing loop partitioning. The program generator inserts code before each parallel loop to maintain the coherence of shared data (if necessary, see section 3.6) and inserts a barrier after each parallel loop.

The program generator also produces code to allocate distributed arrays according to the data decomposition directives provided in the input program. This code allocates a region in the local memory of each processor to hold each array partition that resides in this local memory. This allocation results in an array which is physically distributed in shared memory, but is accessed in the output SPMD program as if it is allocated in contiguous form. Array references in the data parallel program are converted into equivalent array references to the distributed array in the SPMD program. We refer to this conversion as *reference translation*. The program generator takes advantage of the single address space and uses special array allocation schemes that reduce the run-time overhead of reference translation.

3.5. *Reference translation*

In this section we describe the techniques used to perform reference translation. These techniques are presented using 2-dimensional arrays, and assuming row-major storage order of arrays in memory. The techniques are equally applicable to higher dimension arrays and to column-major storage order.

Reference translation for an array is performed using a number of allocation schemes for the partitions of the array and using one or more array *descriptors*. A descriptor is an array of pointers in which each pointer points to a region of an array partition. The pointers are initialized by the program generator in such a way to allow the array to be accessed as if it is contiguous, when it is physically distributed in memory. The initialization of the pointers in a descriptor is referred to as *region binding*. Array references in the lhs of statements are always writes to local array elements, while array references in the rhs of statements are reads that may access local or non-local array elements. The translation of lhs array references is called *local reference translation*; the translation of rhs array references is called *global reference translation*.

3.5.1. *Local reference translation* Local reference translation is performed using *local descriptors*. There is one local descriptor for each distributed array in each processor; the dimensionality of the descriptor is one less than the dimensionality of the array. The local descriptor of an array is private to a processor and is used only to access the processor's own portion of the array.

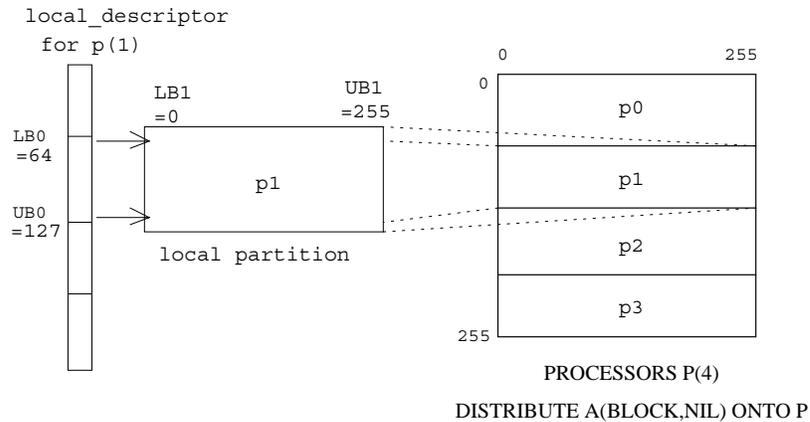


Figure 5. Region binding for (BLOCK, NIL) decomposition.

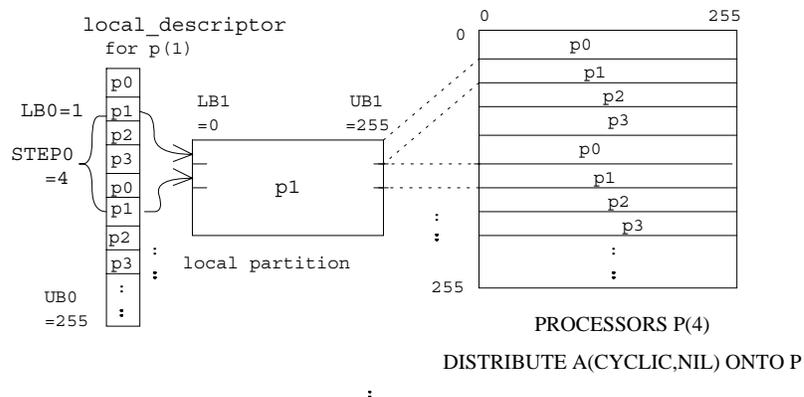


Figure 6. Region binding for (CYCLIC, NIL) decomposition.

The binding of each local descriptor for an array depends on the partitioning attribute assigned the last dimension of the array. If the attribute assigned to the last dimension of the array is NIL, a contiguous local region is allocated in the memory of each processor to hold the rows of each local partition. The binding of the local descriptor is done by assigning the addresses of the rows of the partition to corresponding rows of the local descriptor. This is shown in Figures 5 and 6 for the (BLOCK, NIL) and (CYCLIC, NIL) decompositions respectively.

If the partitioning attribute assigned to the last dimension of an array is BLOCK, a contiguous region of memory is allocated to hold the row segments of the local array partition. Region binding is done by assigning the address of each row of the array partition to the corresponding row in the local descriptor, with an offset adjustment. The offset added to each row pointer corresponds to the starting posi-

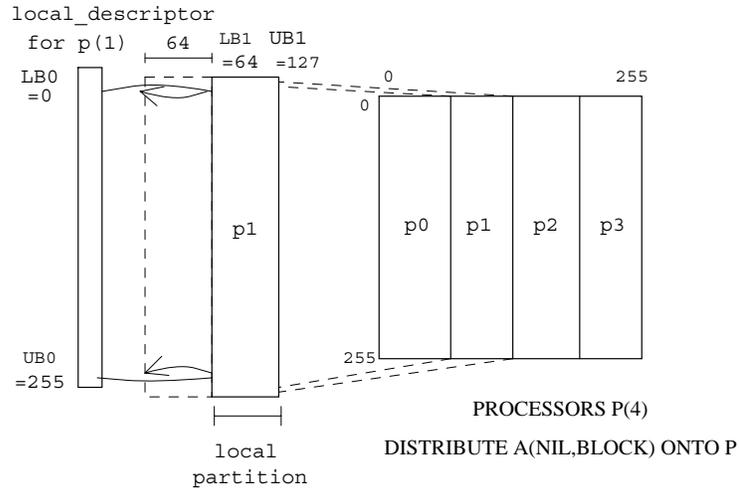


Figure 7. Region binding for (NIL, BLOCK) decomposition.

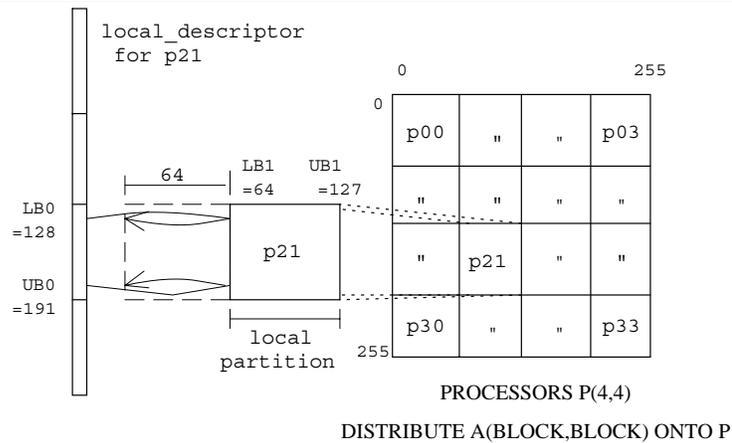


Figure 8. Region binding for (BLOCK, BLOCK) decomposition.

tion of left column boundary of the local array partition. This is illustrated for the (NIL, BLOCK) and (BLOCK, BLOCK) decompositions in Figures 7 and 8 respectively.

If the partitioning attribute assigned to the last dimension of any array is CYCLIC, the whole row is allocated and bound to the local descriptor. The processor uses only a subset of this region, which corresponds to the elements it owns; the remaining elements are not used. Hence, this allocation scheme results in memory overhead. Compared to the previous two cases in which only the amount of storage needed by a processor is allocated, this scheme for region binding is less space-efficient. In section 3.7, a number of optimizations will be presented to eliminate

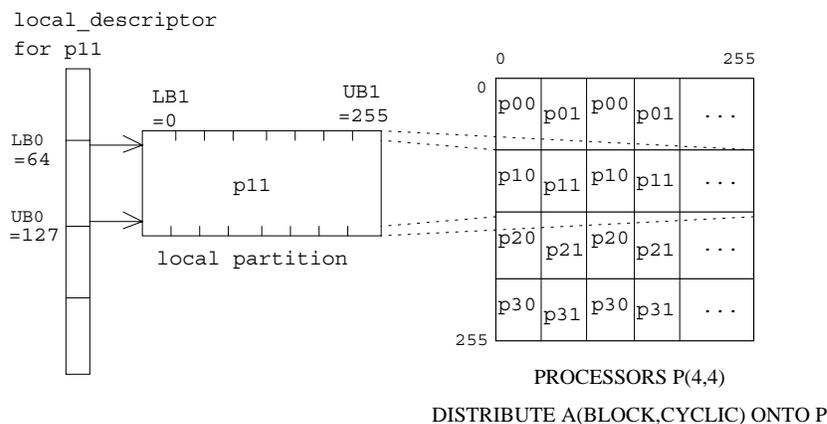


Figure 9. Region binding for (BLOCK, CYCLIC) decomposition.

this memory overhead. Region binding for the (BLOCK, CYCLIC) decomposition is shown in Figure 9.

3.5.2. Global reference translation Global reference translation is used to translate references to arrays in the rhs of statements, which are read references to arrays elements that are both local and non-local to a processor. *Global descriptors* are used to perform global address translation. Similar to a local descriptor, the dimensionality of a global descriptor is one less than the dimensionality of its corresponding array. However, unlike local descriptors, global descriptors are shared variables; there is only one global descriptor for each distributed array, and all processors share this descriptor.

There are two methods for binding the global descriptor of any array to the various partitions of the distributed array. *Remote region binding* is used when the NIL attribute is assigned to the last dimension of an array. In this case, each row pointer is assigned to the corresponding row's location in the global descriptor. This is depicted in Figure 10 for (BLOCK, NIL) decomposition. *Image binding* is used when the partitioning attribute assigned to the last dimension of the array is not NIL. An *image array* is an array of pointers of the same size and dimensions as the data array. Each element in image array is a pointer to the corresponding element of the distributed array. The image array is used to access the elements of the distributed array through pointer dereferencing. The use of image binding is depicted in Figure 11 for (BLOCK, BLOCK) decomposition.

The use of the image array to perform global reference translation results in memory overhead since an additional image array is used for each data array. In section 3.7 a number of optimizations to eliminate this overhead will be described.

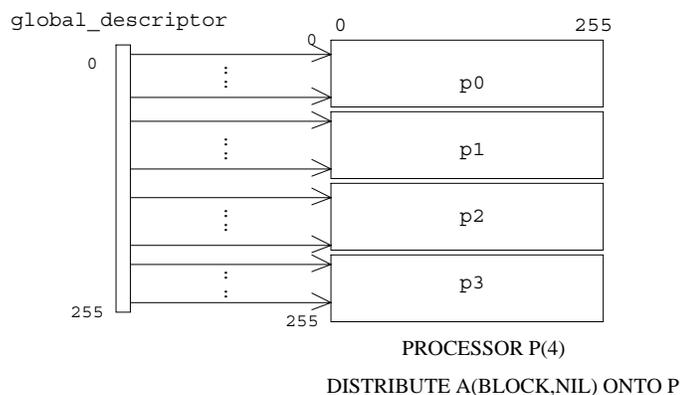


Figure 10. Remote region binding for (BLOCK, NIL) decomposition.

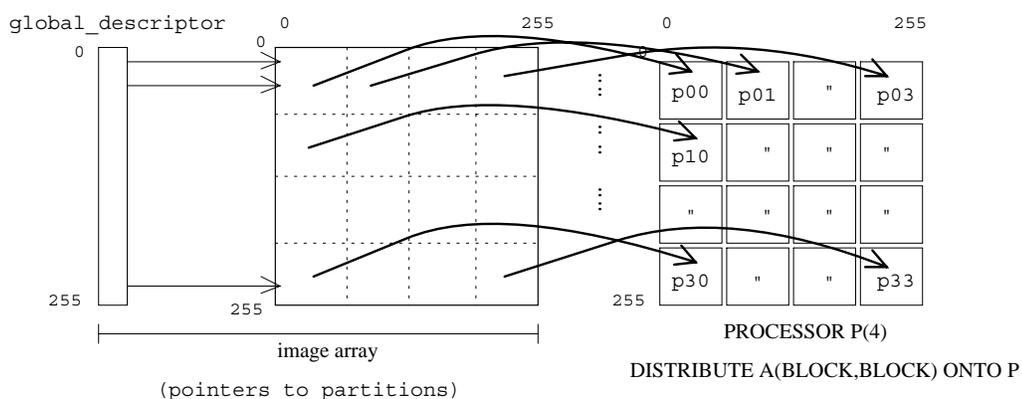


Figure 11. Image binding for (BLOCK, BLOCK) decomposition.

3.6. Shadow regions

In this section, we describe a scheme to reduce false sharing in distributed arrays. This scheme is orthogonal to the data distribution schemes, in that it can be used with any data decomposition scheme.

In this scheme, the compiler uses *shadow regions* to maintain cache consistency. Shadow regions are local memory storage used to hold non-local data in the local memory of a processor during the execution of a parallel loop. A processor accesses non-local data by accessing the local shadow region. The program generator inserts code to maintain memory consistency of data between local shadow regions and non-local array partitions.

Program analysis is used to determine interprocessor communication for each distributed array on each processor. This allows the compiler to determine the

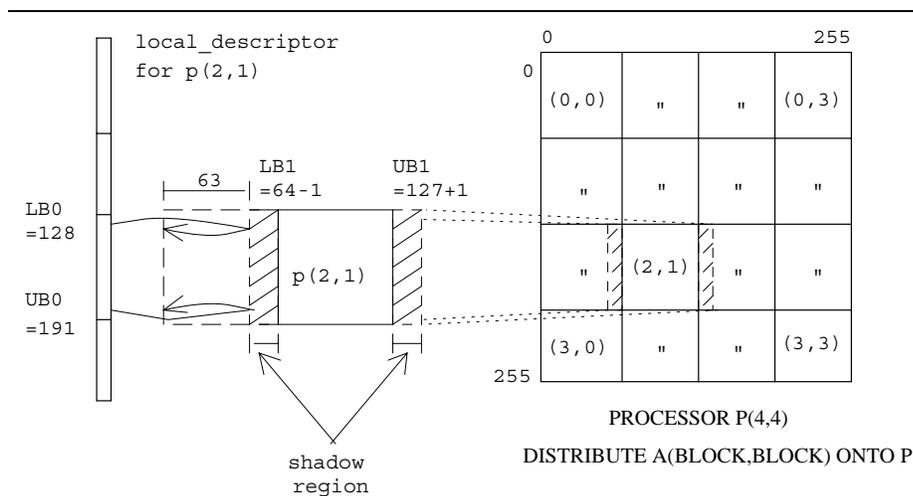


Figure 12. Shadow region binding for (BLOCK, BLOCK) decomposition.

size and location of local buffers required to hold the non-local data. The compiler allocates a memory region in each processors that is large enough to hold both the local array partition as well as the shadow region. The compiler generates code such that the remote data is accessed correctly in the program by copying data from data partition into the shadow region before each parallel loop.

Once the shadow region is created, it is bound to a local descriptor variable. This binding process is called *shadow region binding*. The process of binding a shadow region is similar to that of binding a local partition. Figure 12 depicts shadow region binding for (BLOCK, BLOCK) decomposition.

The compiler generates prefetching code to copy non-local data into the shadow region before each parallel loop. This code employs a two-step copying process that leaves both local and non-local data partitions cacheable. The processor containing the non-local data first copies the data into a temporary buffer. The processor accessing the non-local data then copies the temporary buffer into the shadow region, which leaves only the buffer uncacheable. This copying process also maintains data coherence between the non-local memory and the local shadow memory through the uncacheability of the buffer. The non-local processor invalidates the appropriate cache lines out of its cache as it writes the data to the temporary buffer and the local processor reads the data from memory.

The use of shadow regions has two main advantages. First, a shadow region provides faster access to non-local data. Array references in the rhs of a loop statement are translated into local shadow region accesses. A non-local array element is copied once into the local shadow region of a processor, and is then read by the processor from the shadow region. Hence, when there is reuse of non-local array data in a parallel loop, the cost of non-local memory access is incurred only once rather than every time the non-local element is accessed. Second, since the shadow

and the local partition are merged in one local region, there is no longer a distinction between local and non-local accesses, which simplifies generated parallel code by avoiding the possible use of image arrays.

However, shadow regions result in memory overhead since they are implemented as extensions to local partitions. This overhead is small in parallel applications with high spatial locality in which a processor accesses “border” data from neighboring processors. The overhead can be large in applications which exhibit little spatial locality making non-local accesses sparsely distributed over the array; shadow region elements which are not used to store data used by the processor are wasted.

3.7. Reduction of memory overhead

Memory overhead results when an image array is used in global reference translation. This overhead not only increases the amount of memory used by a program, but also degrades program performance. The use of the image array increases data access time because of pointer dereferencing. Also, image array elements can eject data array elements from the local cache either due to limited cache capacity or due to conflicts, which, in turn, increases the number of cache misses incurred by the program.

In this section, we present two optimizations to reduce memory overhead by eliminating the need for the imaging array: dimension swapping and loop splitting.

3.7.1. Dimension swapping Dimension swapping is an optimization that allows global reference translation to be performed using remote region binding instead of image binding, hence, eliminating the overhead of the image array altogether. Consider an n -dimensional array $A[s_0][s_1] \dots [s_{n-1}]$ with the partitioning scheme $P(p_0, p_1, \dots, p_{n-1})$, where s_0, \dots, s_{n-1} are the sizes of the respective dimensions of A , and p_0, \dots, p_{n-1} are the partitioning attributes of the respective dimensions of A . Dimension swapping can be used when the last partitioning attribute p_{n-1} is not NIL and one of the other attributes is NIL. Let the dimension to which the NIL attribute is assigned be p_d ; that is, $p_d = \text{NIL}$, $0 \leq d \leq n-2$. The array A is transformed into an array $A'[s'_0][s'_1] \dots [s'_d] \dots [s'_{n-2}][s'_{n-1}]$ with the new partitioning scheme $P'(p'_0, p'_1, \dots, p'_d, \dots, p'_{n-2}, p'_{n-1})$ such that

$$\begin{aligned} s'_i &= s_i; & 0 \leq i \leq n-2, i \neq d \\ s'_d &= s_{n-1}; \\ s'_{n-1} &= s_d; \end{aligned}$$

$$\begin{aligned} p'_i &= p_i; & 0 \leq i \leq n-2, i \neq d \\ p'_d &= p_{n-1}; \\ p'_{n-1} &= p_d. \end{aligned}$$

That is, the new array A' is expressed as $A'[s_0][s_1] \dots [s_{n-1}] \dots [s_{n-2}][s_d]$, and has the partitioning scheme $P'(p_0, p_1, \dots, p_{n-1}, \dots, p_{n-2}, p_d)$. In the program, a reference to the array A is replaced by a reference to A' in which the subscript expressions in dimensions d and $n - 1$ are swapped. For example, a reference in the form $A(i_0, i_1, \dots, i_d, \dots, i_{n-2}, i_{n-1})$ is replaced by a reference in the form $A'(i_0, i_1, \dots, i_{n-1}, \dots, i_{n-2}, i_d)$.

Dimension swapping improves performance by eliminating the the overhead of the image array. However, it may lead non-unit stride accesses of the array, which can degrade cache performance. The extent to which this optimization improves performance is both application and target-machine dependent.

3.7.2. Loop splitting Loop splitting is an alternative optimization for replacing image array binding with remote region binding. It is used when dimension swapping is not possible. In this optimization, loop iterations that cause rhs array references in a statement to read non-local data are peeled off the loop [24]. The peeling process transforms each loops nest into an equivalent sequence of loop nests such that in each of the resulting loop nest, rhs array references are to one and only one processor. This allows the compiler to replace all references to the image array by direct references to the corresponding regions of the array.

The iterations that need to be peeled off a loop can be easily determined using the inverse of rhs array subscript expressions in the loop. However, it is also necessary to determine the ownership of the remote elements accessed in each iteration by a given reference. The algorithm shown in Figure 13 is used for this purpose. The algorithm takes the set of loop iterations that must be peeled, the set of rhs array references and determines for each array reference which loop iterations cause a reference on each processor. This information is then used to split the loops as described above.

Loop splitting improves performance by eliminating the need for image array binding. However, it degrades performance because of increased loop overhead. The extent to which it benefits performance will depend on loop overhead cost relative to image array cost. Since the number of iterations peeled off a loop depends on the communication patterns in the application, the relative benefit of this optimization is also application and target-machine dependent.

4. Experimental results

A prototype of the compiler described above has been implemented for the Hector Multiprocessor and applied to 4 benchmark kernels: Jacobi, Matrix Multiply, LU factorization, and Simplex. The implementation takes as input programs written in a restricted version of the C programming language [25], and produces C code with Hurricane system calls to create and synchronize parallel threads. In this section, we describe the experimental results obtained from a 16-processor cluster consisting of 4 stations on a single local ring. The results demonstrate the performance improvements that can be obtained by using our compiler on the application bench-

Algorithm OwnershipLoopSplit

Input: peeled iteration set \mathcal{L}

a set of rhs array references \mathcal{R}

a list of processor locations $\mathcal{P} = \text{NIL}$

Output: a set of tuples $\mathcal{S} = \{(\mathcal{X}, \mathcal{M})\}$

where $\mathcal{X} \in \mathcal{L}$ and \mathcal{M} denotes the processor locations

begin

let $R =$ first reference in \mathcal{R}

for each processor location p for \mathcal{R}

find the iteration set $\mathcal{X} \in \mathcal{L}$ such that R causes a reference to the processor p

if $\mathcal{X} \neq \text{NIL}$

let $\mathcal{M} = \mathcal{P}$

add p to \mathcal{M}

if $|\mathcal{R}| == 1$

add $(\mathcal{X}, \mathcal{M})$ to the set \mathcal{S}

else /* $|\mathcal{R}| > 1$ */

let $\mathcal{R}' = \mathcal{R} - R$

call LoopPeel($\mathcal{X}, \mathcal{R}', \mathcal{M}$)

end

Figure 13. The algorithm for determining ownership in loop splitting.

marks. These improvements compare favorably to the improvements obtained using the operating system data management policies, and to the improvements obtained using a set of data management macros called Numacros (described in the next section). The results also show the impact of memory reduction optimizations on performance.

The performance of the compiler-generated parallel kernels on Hector is shown in Figure 14. The performance is measured by the speedup, which is defined as the ratio of the execution time of the sequential kernel to the execution time of the compiler-generated parallel kernel. The data sizes for each kernel used to measure performance is selected so as all data used by the sequential kernel fits in the local memory of one processor module. This is done to ensure that speedup numbers are not inflated due to the increase in memory locality when more than one processor is used; using a large data set would make the sequential program slower due to the need for non-local memory allocation, and consequently non-local memory accesses.

The results generally indicate that the speedup of the compiler-generated parallel codes scales well with the number of processors. The extent to which the performance of each kernel improves, and the data decomposition scheme needed to achieve such improvement varies from one application to another, and depends on the number of processors.

The performance of Jacobi is shown for three input data sizes (i.e., sizes of arrays a and b), 350×700 , 400×400 and 700×350 . The performance of the kernel without the use of shadow regions scales well for the (BLOCK,NIL) scheme, but is

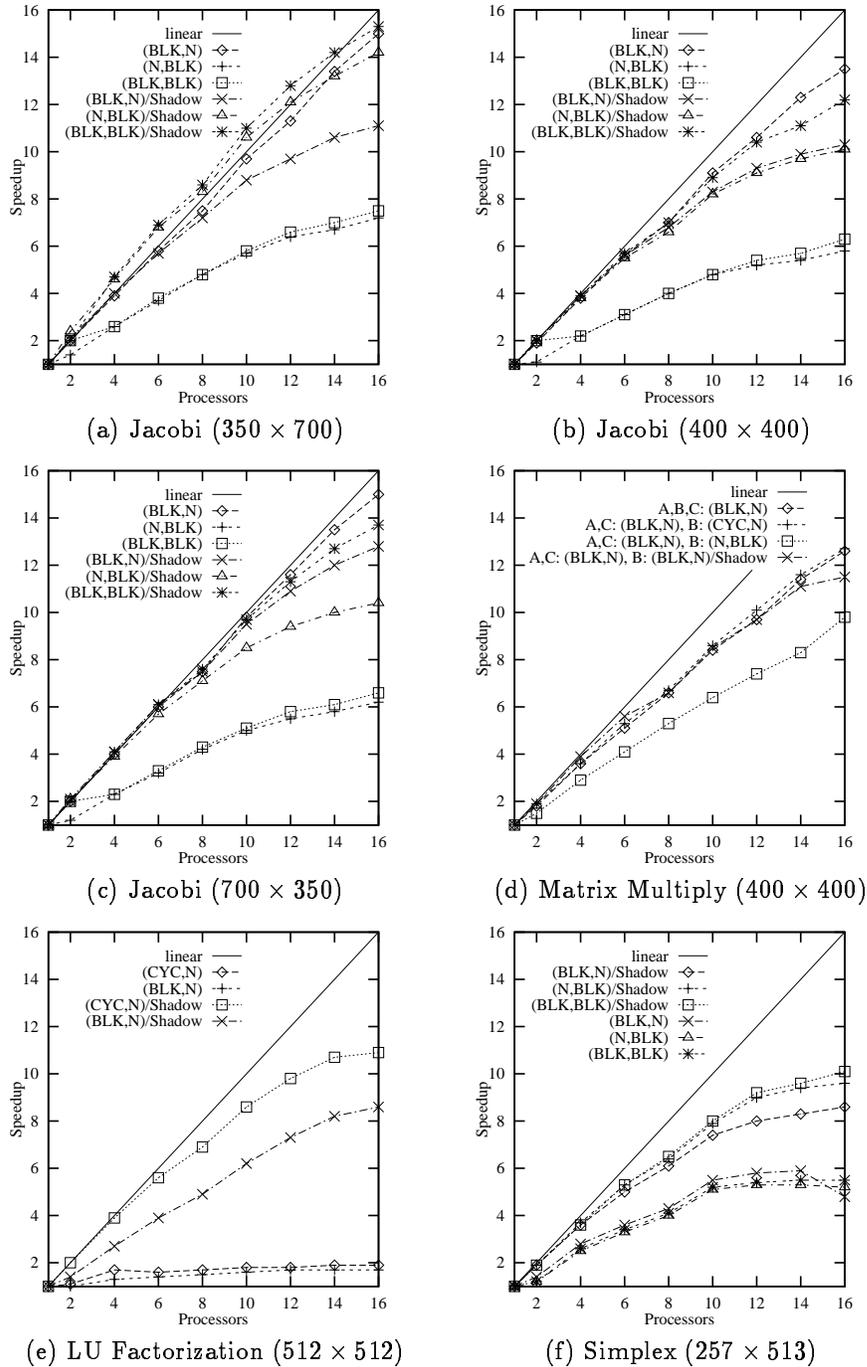


Figure 14. Performance of compiler generated parallel programs.

poor for the (NIL,BLOCK) and (BLOCK,BLOCK) schemes. The poor performance for these two schemes is attributed to the resulting pattern of remote data accesses, which causes uncacheability of data. Since arrays are stored in row-major order, accessing the non-local data in column-major order causes most, if not all, of the pages in non-local array partitions to become uncacheable. In contrast, the (BLOCK,NIL) scheme makes uncacheable only pages containing boundary rows in the partitions, which does not impact on performance in a significant way. In addition, the (NIL,BLOCK) and (BLOCK,BLOCK) schemes require the use of an image array for remote reference translation, which further degrades performance.

The performance of Jacobi with the use of shadow regions is generally scaling for all data decomposition scheme. The performance is largely determined by the amount of remote data communication in the parallel program. In the case of Jacobi, this is determined by the number of data elements shared between pairs of processors. In the case of the (NIL,BLOCK)/Shadow scheme, boundary columns in the arrays *a* and *b* are shared, and the scheme performs better than the (BLOCK,NIL)/Shadow scheme in the case where the column size is smaller than the row size of the array; i.e. 350×700 . Similarly, (BLOCK,NIL)/Shadow scheme performs better than the (NIL,BLOCK)/Shadow scheme when the row size is smaller than the column size, i.e., 700×350 . As expected, the (BLOCK,BLOCK)/Shadow scheme out-performs both in for all data sizes since it always minimizes communication volume.

The performance of the Jacobi kernel with the use of shadow regions is better than its performance without shadow regions when false sharing is a major factor affecting performance. The implementation of shadow regions overcomes false sharing, resulting in better performance. However, when false sharing is not a major factor, the performance of the kernel without the use of shadow regions is better than its performance with shadow regions. This is due to the overhead associated with data copying in the case of shadow regions.

The superlinear behavior exhibited by the Jacobi kernel for the data size of 350×700 is due to the decrease in the size of the data handled per processor as larger number of processors are used, and the increase opportunity for data re-use offered by the column and grid block distributions. The decreasing size of the local data makes it more likely for the data to fit in the local cache.

The performance of the matrix multiply kernel is shown in Figure 14(d). A (BLOCK,NIL) distribution is used for arrays *A* and *C* since it maximizes the locality of these arrays. The array *B*, which is accessed in its entirety by all processors is distributed using a number of schemes. The figure indicates that the schemes, with the exception of (NIL,BLOCK) result in good performance. The poor performance in the case of the (NIL,BLOCK) scheme is due to the need to use the image array to perform global reference translation. The (NIL,CYCLIC) distribution results in similar performance, and is not shown.

The performance of the matrix multiply kernel for the remaining schemes is fairly comparable. The shadow scheme performs slightly better for a small number (≤ 4) of processors when compared to the other schemes, but does not scale well as more processors are used. This is because each processor needs to access all

of array B multiple times (once per i iteration). The size of the shadow region needed is equal to the size of the array B. The communication code generated by the compiler copies B into the shadow region of each processor before the parallel loop. Although the resulting shadow region memory accesses are faster than remote memory accesses, the cost of broadcasting B becomes large when the number of processors increases. In addition, the number of times each processor is required to access B is inversely proportional to the total number of processors used (fewer i iterations). Consequently, as the number of processors increases, the reuse of the shadow data decreases and with it the performance advantage of using the shadow policy.

The performance of the LU kernel is shown in Figure 14(e). The speedup of the kernel using a block distribution is extremely poor because of the load imbalance caused by the triangular iteration space of the parallel loop in the kernel. The cyclic distribution performs better since it improves load balance. We have shown only the row-cyclic distribution since column-cyclic distributions cause false sharing, which degrades performance, as seen in Jacobi above. The performance using a shadow region with the row-cyclic distribution is better than that using the operating system cache coherence policy with the same distribution. Similar to Jacobi, this is because of false sharing.

However, the performance of the LU kernel using shadow regions does not scale well with increasing number of processors. This is mainly due to the size of the shadow region generated by the compiler, which is equal to the size of the portion of the array a being updated in the Gauss-Jordan elimination step. This large shadow region is copied once in each k iteration of the kernel; i.e, $n-1$ times, which affects performance negatively.

The performance of the Simplex kernel is shown in Figure 14(f). The use of shadow regions results in better performance than the operating system policies. This is attributed to the false sharing introduced by the operating system page-based cache coherence policy.

In Figure 15, the performance of the parallel programs generated by our compiler is compared to the performance of the parallel programs generated by Numacros, as well as to the performance of manually parallelized programs that use the default Hurricane memory management policies. The figure shows the results for the data size of 512×512 , which is the only size for which Numacros results are available. For each kernel, the figure compares the performance obtained by our compiler to the best performance obtained by Numacros and to the best performance using the Hurricane policies. Numacros results are not available for the Simplex kernel. The figure shows that the performance of our compiler generated programs are significantly better than the performance of the Hurricane-based programs. The figure also shows that the best performance obtained by our compiler is better than that of Numacros, except for LU. In this case, Numacros performs better because of its use of optimizations that reduce the number of the integer division operations. However, the Numacros performance is obtained manually matching loop and data partitioning; in the case of our compiler, loop partitioning is automatically determined by the compiler.

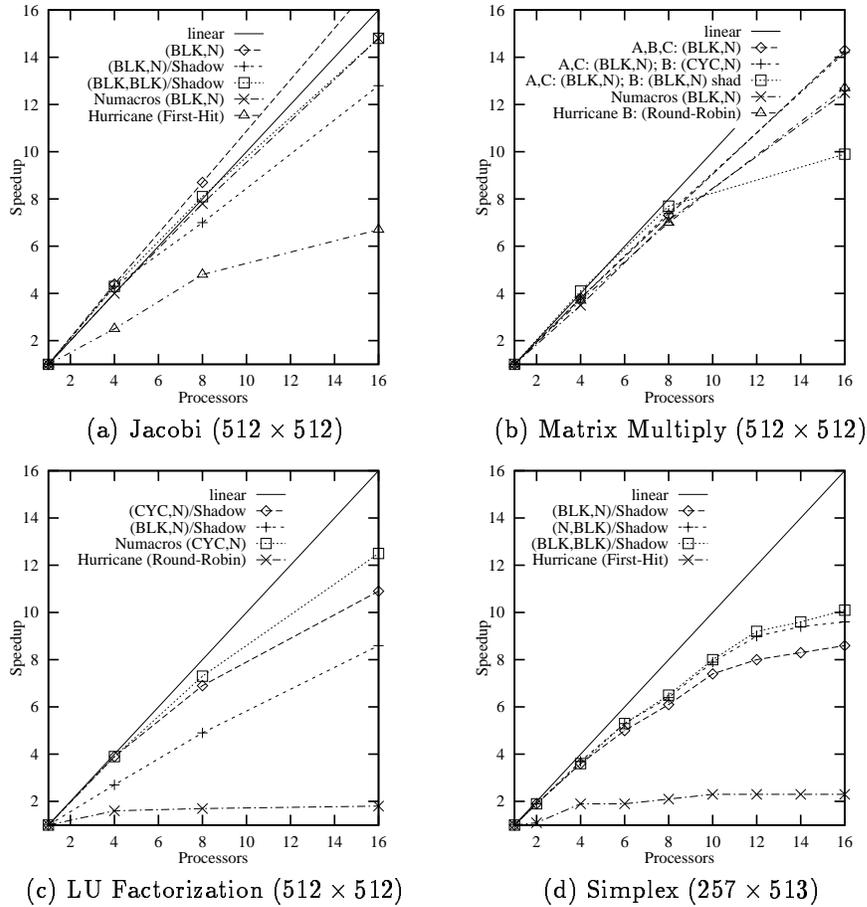


Figure 15. Performance comparisons.

The impact of memory overhead reduction optimizations on the performance of the Jacobi kernel can be seen in Figure 16. The figure shows the speedup of the Jacobi kernel using a (NIL,BLOCK) decomposition scheme for both arrays *a* and *b*, which requires the use of an image array for global reference translation. The figure shows the speedup of the same kernel after each of the two optimizations described in section 3.7 is applied. Dimension swapping eliminates the need for the image array without introducing any overhead, and hence, results in extremely good performance. Loop splitting can be used to break the first loop of the kernel into three loops such that remote region binding can be used in each loop, eliminating the need for the image array. This results in good performance, but the increase in loop overhead makes the optimization less effective than dimension swapping. Hence, the above results suggest that the compiler's strategy for applying the memory

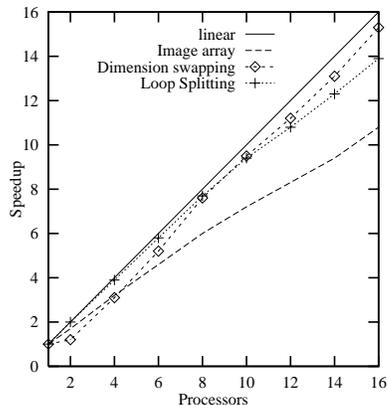


Figure 16. Impact of memory overhead reduction optimizations on performance.

overhead optimizations should be to first attempt dimension swapping, and if not possible, apply loop splitting.

5. Related work

A number of research and production compilers have been built to compile a data-parallel program into an SPMD parallel program for distributed memory multiprocessors. They include Fortran D [9], High-Performance Fortran [10], SuperB [26], and Vienna Fortran [6]. These compiler share the common approach of using programmer-supplied directives for specifying array decomposition schemes. Our compiler shares the same approach, but targets NUMA shared memory multiprocessors, hence, extending the applicability of the approach beyond distributed memory multiprocessors. Our compiler is unique in the way it performs address translation taking advantage of the existing shared address space.

A number of research compilers have also been built to tackle the problem of automatically deriving array decompositions in programs [2, 11, 21]. This problem is orthogonal to the one considered in this paper. The aim of such compilers is to determine *what* array decompositions are suitable for a given program. Our work describes *how* to efficiently support array decompositions on NUMA multiprocessors, and demonstrates their advantages over operating system policies.

A number of approaches have been proposed to manage shared data on NUMA multiprocessors. They include operating system policies [20, 4, 5, 22], run-time libraries [17, 18] and user-level macros [15].

Stumm, Unrau and Krieger [20] describe the design of the Hurricane operating system in which clusters of loosely-coupled micro-kernels are used instead of a single monolithic kernel to enhance locality. They show that this clustering approach leads to better locality. Our experiments have been conducted using Hurricane, and we

show that additional improvements in application performance can be attained using our proposed methods.

LaRowe et al. [12, 13] experiment with a number of operating system page placement policies, including page migration, on a BBN Butterfly. Their experiments indicate that a general-purpose page placement policy is inferior to application-oriented data placement, mainly due to false sharing and suboptimal placement of data. Furthermore, their experiments indicate that while their general-purpose policy results in adequate performance, further improvements can be attained by tuning the policy to the data access pattern of a particular application. Our proposed scheme enhances locality using the compiler's knowledge of data access patterns and using the supplied directives, achieving application-oriented data placement.

Chandra et al. [4, 5] explore process scheduling strategies and page migration policies on Unix-based multiprocessor systems. They show that their affinity-based scheduling policy outperforms Unix scheduling policies for workloads consisting of sequential jobs. However, for workloads consisting of parallel applications, the best policy depends on application characteristics. They also show that a dynamic page migration policy based on TLB misses is viable. More recently, Verghese et al. [22] extend Chandra's work and describe support for improving locality in the Hive operating system. They allow dynamic page migration and replication based on cache misses instead of TLB misses. They show that both page migration and replication are necessary and that using both can improve execution time by as much as 29% over a First-Hit policy. Both approaches focus on improving the performance of multiprogrammed workloads, and manage data at the granularity of the page, which is too coarse to capture the fine-grain sharing in parallel applications. In contrast, our work focuses on improving locality for a single application and allows data to be managed at a granularity smaller than a page.

Sandhu, Gamsa and Zhou [17, 18] also address the issue of granularity and propose *shared regions* as a paradigm for data management on NUMA shared memory multiprocessors. Programmers annotate their parallel code by specifying the regions accessed by a thread as part of the synchronization primitives in the program. A run-time system migrates and replicates data across the physical memories in the system to enhance locality. Shared regions can lead to good performance, but require programmers to first parallelize and then annotate applications. In contrast, our approach relies only on an abstract specification of data distribution and parallelism; the compiler automatically generates a parallel program.

Li and Sevcik [15] describe a set of C language macros that implement data distribution and loop partitioning for improved performance on NUMA multiprocessors. Their approach is similar to ours in that they use data decompositions to achieve locality. However, they require programmers to manually insert loop partitioning directives in addition to data partitioning directives, and they use array allocation strategies that lead to additional index calculation. They propose compiler optimizations to reduce the impact of these calculations. In addition, they ignore cache performance by flushing cache contents after each parallel loop. In contrast, we provide shadow regions as a mechanism for improving cache performance; data partitions remain cacheable, only temporary buffers are flushed out of the cache.

6. Conclusions

Management of program data to reduce false sharing and improve memory locality is important for scaling performance of applications on NUMA shared memory multiprocessors. However, restructuring programs to manage data is tedious and error-prone. In this paper, we suggested and demonstrated that data management in array-based scientific applications can be effectively performed using compilers, making it easy for programmers to manage data. We have developed a compiler system that accepts as input a data parallel program in which parallelism is expressed using the FORALL statement, as well as data decomposition schemes for arrays. The compiler generates an SPMD program in which parallelism is expressed using appropriate operating system calls, and in which arrays are partitioned and placed in the physical memory modules of the multiprocessor. The compiler takes advantage of the underlying shared memory architecture to efficiently implement array partitioning and distribution. In addition, the compiler provides a cache-coherence scheme based on shadow regions to avoid false sharing. Experimental results obtained from the Hector multiprocessor indicated that the compiler-generated parallel programs have better performance compared to parallel programs which rely on operating systems polices for data management.

The techniques presented in this paper are applicable to other multiprocessor platforms as well. The performance of applications on shared memory multiprocessors with hardware support for cache coherency can be improved by reducing false sharing that occurs at the cache-line level and by enhancing memory locality. The techniques are also applicable on networks of workstations that employ virtual-shared memory to provide programmers with a single coherent view of memory. On such systems, our techniques can be expected to significantly improve performance as interprocessor communications is performed using costly message passing.

The study indicates that the selection of a data decomposition scheme is critical for scaling performance. However, the process of arriving at the most suitable data distribution can be tedious in large programs. Hence, one extension to this work is to enhance the compiler with a static performance estimator that provides programmers with feedback on the quality of selected distributions. Another extension is to address the support for dynamic data decompositions to accommodate changing data access patterns during run-time.

Acknowledgments

This research is supported by grants from the Natural Sciences and Engineering Council of Canada and from the Information Technology Research Center of Ontario.

References

1. T. Abdelrahman and T. Wong. Distributed array data management on NUMA multiprocessors. In *Proceedings of Scalable High Performance Computing Conference*, pages 551–559,

- 1994.
2. J. Anderson and M. Lam. Global optimizations for parallelism and locality on scalable parallel machines. In *Proceedings of ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 112–125, 1993.
 3. V. Balasundaram. A mechanism for keeping useful internal information in parallel programming tools: The data access descriptor. *Journal of Parallel and Distributed Computing*, 9(6):154–170, 1990.
 4. R. Chandra. *The COOL parallel programming language: design, implementation and performance*. Ph.D. Thesis, Department of Computer Science, Stanford University, 1995.
 5. R. Chandra, S. Devine, B. Verghese, A. Gupta and M. Rosenblum. Scheduling and page migration for multiprocessor compute servers. In *Proceedings of the 6th Int'l Conference on Architectural Support for Programming Languages and Operating Systems*, pages 12–24, 1994.
 6. B. Chapman, P. Mehrotra and H. Zima. Vienna Fortran— a language extension for distributed memory multiprocessors. Technical report 91–72, ICASE, 1991.
 7. Convex Computer Corporation. *Convex Exemplar System Overview*. Document No. 080-002293-000, Richardson, TX, USA, 1994.
 8. Cray Research. *The T3D massively parallel processor system*. Cray Research, 1993.
 9. S. Hiranandani, K. Kennedy and C. Tseng. Compiling Fortran D. *Communications of the ACM*, 35(8):66–80, 1992.
 10. C. Koelbel, D. Loveman, R. Schreiber, G. Steele, Jr., and M. Zosel. *The High Performance Fortran Handbook*. The MIT Press, Cambridge, MA 1994.
 11. U. Kremer. *Automatic data layout for distributed memory machines*. Ph.D. Thesis, Department of Computer Science, Rice University, 1995.
 12. R. LaRowe, C. Ellis, and L. Kaplan. The robustness of NUMA memory management. In *Proceedings of the Symposium on Operating System Principles*, pages 137–151, 1991.
 13. R. LaRowe, J. Wilkes and C. Ellis. Exploiting operating system support for dynamic page placement on NUMA shared memory multiprocessors. In *Proceedings of the 3rd ACM Symposium on Principles and Practice of Parallel Programming*, pages 122–132, 1991.
 14. D. Lenoski, J. Laudon, K. Gharachorloo, W. Weber, A. Gupta, J. Hennessy, M. Horowitz and M. Lam. The Stanford Dash multiprocessor. *IEEE Computer*, 25(3):63–79, 1992.
 15. H. Li and K. Sevcik. NUMACROS: Data Parallel Programming on NUMA Multiprocessors. In *Proceedings of 4th Symposium on Experiences with Distributed and Multiprocessor Systems (SEDMS IV)*, pages 247–263, 1993.
 16. D. Palmero and P. Banerjee. Automatic selection of dynamic data partitioning schemes for distributed-memory multicomputers. In *Workshop on Languages and Compilers for Parallel Computing*, pages 392–406, 1995.
 17. H. Sandu. *Shared regions: a strategy for efficient cache management in shared-memory multiprocessors*. Ph.D. Thesis, Department of Computer Science, University of Toronto, 1995.
 18. H. Sandhu, B. Gamsa and S. Zhou. The shared region approach to software cache coherence on multiprocessors. In *Proceedings of the 4th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 229–238, 1993.
 19. D. Solow, *Linear programming: An introduction to finite improvement algorithms*, North-Holland, New York, 1984.
 20. M. Stumm, R. Unrau and O. Krieger. Clustering micro-kernels for scalability. In *USENIX Workshop on Micro-Kernels*, pages 285–303, 1992.
 21. S. Tandri and T. Abdelrahman. Automatic partitioning of data and computations on scalable shared memory multiprocessors. In *Proceedings of the Int'l Conference on Parallel Processing*, pages 64–73, 1997.
 22. B. Verghese, S. Devine, A. Gupta and M. Rosenblum. OS support for improving data locality on CC-NUMA compute servers. In *Proceedings of the 7th Int'l Conference on Architectural Support for Programming Languages and Operating Systems*, pages 279–289, 1996.
 23. Z. Vranesic, M. Stumm, R. White and D. Lewis. The Hector multiprocessor. *IEEE Computer*, 24(1):72–79, 1991.
 24. M. Wolfe. *High performance compilers for parallel computing*. Addison-Wesley Publishing Company, Redwood City, CA, 1996.

25. T. Wong. *Data partitioning based compiling techniques for NUMA shared memory multiprocessors*. M.A.Sc. Thesis, Department of Electrical and Computer Engineering, University of Toronto, 1994.
26. H. Zima, H. Bast and M. Gerndt. SUPERB: A tool for semi-automatic MIMD/SIMD parallelization. *Parallel Computing*, 6(1):1-18, 1988.