

## Generating Local Addresses and Communication Sets for Data-Parallel Programs

Siddhartha Chatterjee  
John R. Gilbert  
Fred J. E. Long  
Robert Schreiber  
Shang-Hua Teng

(NASA-CR-194605) GENERATING LOCAL  
ADDRESSES AND COMMUNICATION SETS  
FOR DATA-PARALLEL PROGRAMS  
(Research Inst. for Advanced  
Computer Science) 11 p

N94-15796

Unclas

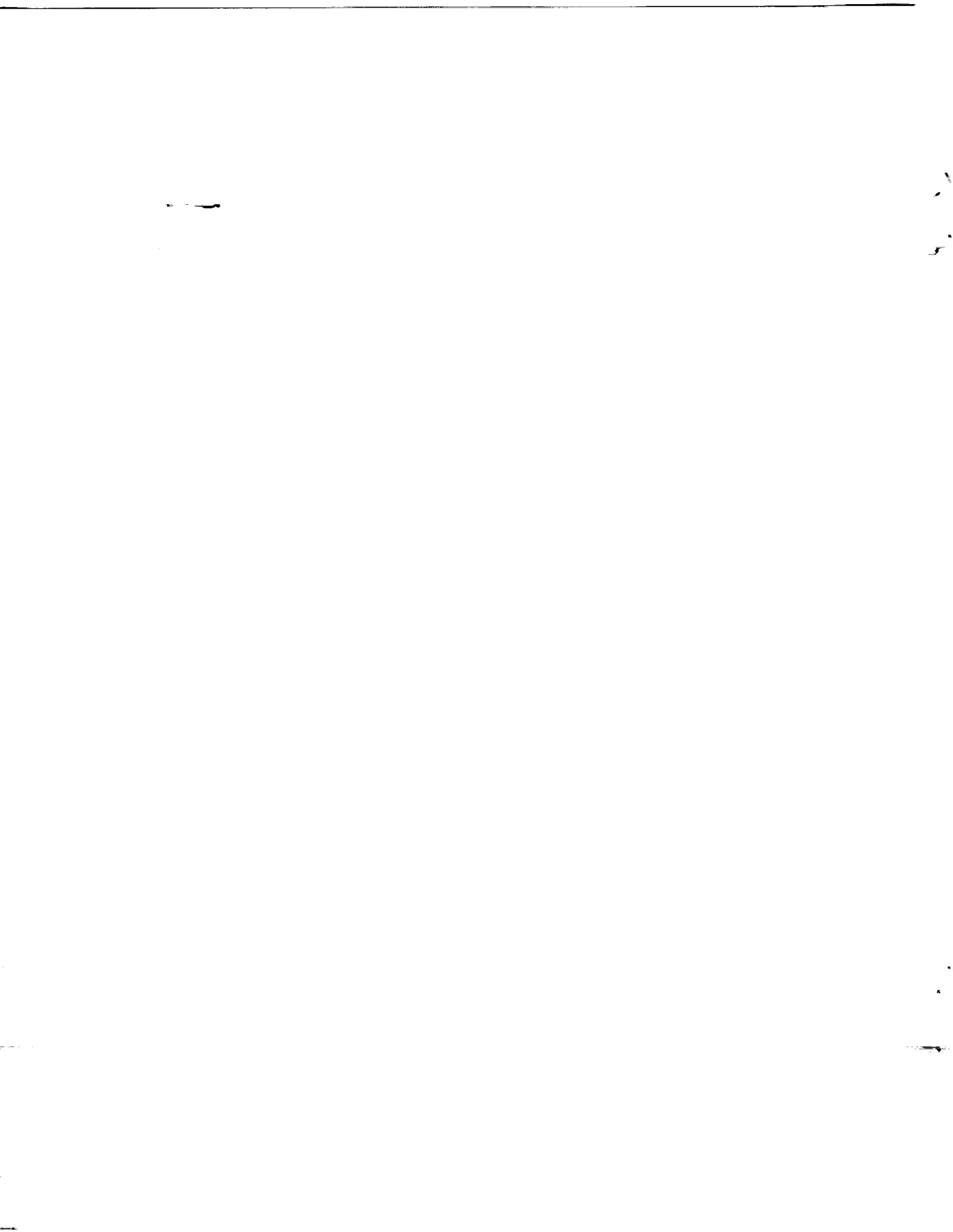
Q3

~~2~~/62 0191135

RIACS Technical Report 93.03

April 1993

To appear in the *Proceedings of the Fourth ACM SIGPLAN Symposium on Principles and Practice of  
Parallel Programming, San Diego, CA, May 1993.*  
Submitted to *Journal of Parallel and Distributed Computing.*



# Generating Local Addresses and Communication Sets for Data-Parallel Programs

Siddhartha Chatterjee  
John R. Gilbert  
Fred J. E. Long  
Robert Schreiber  
Shang-Hua Teng

The Research Institute of Advanced Computer Science is operated by Universities Space Research Association, The American City Building, Suite 311, Columbia, MD 21044, (301) 730-2656

---

Work reported herein was supported by NASA via Cooperative Agreement NCC 2-387 between NASA and the Universities Space Research Association (USRA). Work was performed at the Research Institute for Advanced Computer Science (RIACS), NASA Ames Research Center, Moffett Field, CA 94035-1000.



# Generating Local Addresses and Communication Sets for Data-Parallel Programs

Siddhartha Chatterjee\*   John R. Gilbert†   Fred J. E. Long‡   Robert Schreiber\*   Shang-Hua Teng§

## Abstract

Generating local addresses and communication sets is an important issue in distributed-memory implementations of data-parallel languages such as High Performance Fortran. We show that for an array  $A$  affinely aligned to a *template* that is distributed across  $p$  processors with a *cyclic*( $k$ ) distribution, and a computation involving the regular section  $A(\ell : h : s)$ , the local memory access sequence for any processor is characterized by a finite state machine of at most  $k$  states. We present fast algorithms for computing the essential information about these state machines, and extend the framework to handle multidimensional arrays. We also show how to generate communication sets using the state machine approach. Performance results show that this solution requires very little runtime overhead and acceptable preprocessing time.

## 1 Introduction

Languages such as Fortran D [4] and High Performance Fortran [5] allow the programmer to annotate programs with *alignment* and *distribution* statements that describe how to map arrays to processors in a distributed-memory environment. Both languages introduce an auxiliary Cartesian grid called a *template*; arrays are aligned to templates, and templates are distributed across the processors. The alignment of an array to a template and the distribution of the template determine the final *mapping* of the array. The compiler must partition the arrays and produce “node code” for execution on each processor of a distributed-memory parallel computer.

Consider an array  $A$  of size  $n$  aligned to a template such that array element  $A(i)$  is aligned to template cell  $ai + b$ . The template is distributed across  $p$  processors using a *cyclic*( $k$ ) distribution, in

\*Research Institute for Advanced Computer Science, Mail Stop T045-1, NASA Ames Research Center, Moffett Field, CA 94035-1000 (sc@riacs.edu, schreiber@riacs.edu). The work of these authors was supported by the NAS Systems Division via Cooperative Agreement NCC 2-387 between NASA and the Universities Space Research Association (USRA).

†Xerox Palo Alto Research Center, 3333 Coyote Hill Road, Palo Alto, CA 94304-1314 (gilbert@parc.xerox.com). Copyright © 1992 by Xerox Corporation. All rights reserved.

‡Department of Computer and Information Science, University of California, Santa Cruz, CA 95064 (flong@cse.ucsc.edu). This work was performed while the author was a summer student at RIACS.

§Department of Mathematics, Massachusetts Institute of Technology, Cambridge, MA 02139 (steng@theory.lcs.mit.edu). This work was performed while the author was a postdoctoral scientist at Xerox PARC.

To appear in the Proceedings of the Fourth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, San Diego, CA, May 1993.

$A$	a distributed array
$n$	size of array $A$
$a$	stride of alignment of $A$ to the template
$b$	offset of alignment of $A$ to the template
$p$	number of processors to which the template is distributed
$k$	block size of distribution of the template
$\ell$	lower bound of regular section of $A$
$h$	upper bound of regular section of $A$
$s$	stride of regular section of $A$
$m$	processor number

Table 1: Symbols used in this paper.

which template cell  $i$  is held by processor  $(i \text{ div } k) \bmod p$ . (We number array elements, template cells, and processors starting from zero.) Thus the array  $A$  is split into  $p$  local arrays residing in the local memories of the processors. Consider a computation involving the *regular section*  $A(\ell : h : s)$ , i.e., the set of array elements  $\{A(\ell + js) : 0 \leq j \leq (h - \ell)/s\}$ , where  $s > 0$ . (Table 1 summarizes the notation.) In this paper we consider the following problems:

- What is the sequence of local memory addresses that a given processor  $m$  must access while performing its share of the computation?
- What is the set of messages that a given processor  $m$  must send while performing its share of the computation?

The *cyclic*( $k$ ) distributions generalize the familiar *block* and *cyclic* distributions (which are *cyclic*( $\lceil n/p \rceil$ ) and *cyclic*(1) respectively). Dongarra *et al.* [3] have shown these distributions to be essential for efficient dense matrix algorithms on distributed-memory machines. High Performance Fortran includes directives for multidimensional *cyclic*( $k$ ) distribution of templates.

Many special cases of these problems have been solved, but the general solution does not appear in the literature or (to the best of our knowledge) in any implementation. Koelbel and Mehrotra [8] treat special cases where the array is mapped using a *block* or a *cyclic* mapping, and the regular section has unit stride. Koelbel [7] handles non-unit-stride regular sections with *block* or *cyclic* mappings. MacDonald *et al.* [9] discuss the problem in the context of Cray Research's MPP Fortran, and conclude that a general solution requires unacceptable address computation overhead in inner loops. They therefore restrict block sizes and number of processors to powers of two, allowing the use of bit manipulation in address computation. Experimental implementation of our solution shows that address computation and interprocessor communication can be performed efficiently without these restrictions.

The paper is organized as follows. Section 2 solves the problem for a one-level mapping ( $a = 1, b = 0$ ). Section 3 reduces the problem for two-level mappings to combining the solutions from two subproblems involving one-level mappings. Section 4 extends the framework to handle multidimensional arrays, and Section 5 extends it to handle generation of communication sets. Section 6 discusses performance results, and Section 7 presents conclusions.

## 2 One-level mappings

We first consider the problem where a one-dimensional array  $A$  is aligned to a template with  $a = 1$  and  $b = 0$ . This is called a *one-level mapping*. In this case, the array and the template are in one-to-one correspondence, and we specify the problem by the pair  $(p, k)$  describing the distribution of the template across the processors. As shown in Figure 1(a), we visualize the layout of the array in the processor memories as courses of array elements. Three functions  $\mathcal{P}$ ,  $\mathcal{C}$ , and  $\mathcal{O}$  describe the location of array element  $A(i)$ .  $\mathcal{P}(i; p, k)$  is the processor holding  $A(i)$ ,  $\mathcal{C}(i; p, k)$  is the course of blocks within this processor holding  $A(i)$ , and  $\mathcal{O}(i; p, k)$  is the offset of  $A(i)$  within the course. The layout of array elements in the local processor memories is as shown in Figure 1(b).

These functions are defined as follows.

$$\begin{aligned} \mathcal{P}(i; p, k) &= (i \bmod pk) \operatorname{div} k = (i \operatorname{div} k) \bmod p & (1) \\ \mathcal{C}(i; p, k) &= i \operatorname{div} pk = (i \operatorname{div} k) \operatorname{div} p & (2) \\ \mathcal{O}(i; p, k) &= i \bmod k & (3) \end{aligned}$$

Therefore, we have

$$i = pk \cdot \mathcal{C}(i; p, k) + k \cdot \mathcal{P}(i; p, k) + \mathcal{O}(i; p, k).$$

The following function gives the local memory location (with respect to the base of the local array) where element  $A(i)$  is stored within processor  $\mathcal{P}(i; p, k)$ :

$$\mathcal{M}(i; p, k) = k \cdot \mathcal{C}(i; p, k) + \mathcal{O}(i; p, k). \quad (4)$$

We are interested in the sequence of local memory locations that an individual processor, numbered  $m$ , accesses while performing its share of the computation on the regular section  $A(\ell : h : s)$ . We specify the sequence of accesses by its starting location and by the differences  $\Delta\mathcal{M}(i_1, i_2) = \mathcal{M}(i_2; p, k) - \mathcal{M}(i_1; p, k)$ , where  $i_1$  and  $i_2$  are two successive indices of the part of the regular section mapped to processor  $m$ . As the examples in Figure 2 illustrate, the spacing  $\Delta\mathcal{M}$  is not constant.

The key insight is that (for a particular  $p, k$ , and stride  $s$ ) the offset  $\mathcal{O}(i_1; p, k)$  of one element of the regular section determines the offset  $\mathcal{O}(i_2; p, k)$  of the next element of the regular section in the same processor, and also determines the spacing  $\Delta\mathcal{M}(i_1, i_2)$ . Thus we can represent the offset sequence and the  $\Delta\mathcal{M}$  sequence as a simple table indexed by offset. We visualize this table as the transition diagram of a finite state machine (FSM) whose  $k$  states are the offsets. This FSM is particularly simple because it has no input, and each state has only one outgoing transition; in fact, it is just a union of state-disjoint cycles. Thus the sequence of  $\Delta\mathcal{M}$ 's is periodic, with period at most  $k$ .

The contents of the table, or the transitions of the FSM, depend on  $p, k$ , and  $s$ , but not on the section bounds  $\ell$  and  $h$  or the processor number  $m$ . The same FSM describes the offsets for every processor, though each processor sees at most one cycle of the FSM. The first memory location used by the section in processor  $m$  (and its offset, which is the "start state" for that processor's view of the FSM) depends on  $\ell$  and  $m$  as well as  $p, k$ , and  $s$ .

Following equation (4), we can write  $\Delta\mathcal{M}$  in terms of the differences in course and offset:

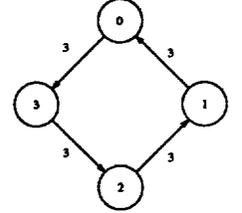
$$\Delta\mathcal{M}(i_1, i_2) = k \cdot \Delta\mathcal{C}(i_1, i_2) + \Delta\mathcal{O}(i_1, i_2), \quad (5)$$

where the definitions of  $\Delta\mathcal{C}$  and  $\Delta\mathcal{O}$  are analogous to that of  $\Delta\mathcal{M}$ .

Here are some examples. We assume that  $\ell = 0$  and that  $h$  is large.

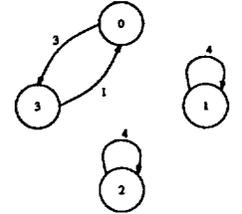
**Example 1**  $k = 4, p = 4, s = 3$ . The distribution is shown in Figure 2(a). The transition table  $T$  is as follows. For convenience, we also tabulate  $\Delta\mathcal{C}$  and  $\Delta\mathcal{O}$ . This shows a simple case where all processors access local memory with constant stride.

STATE	NEXT	$\Delta\mathcal{M}$	$\Delta\mathcal{C}$	$\Delta\mathcal{O}$
0	3	3	0	3
1	0	3	1	-1
2	1	3	1	-1
3	2	3	1	-1



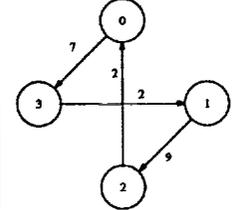
**Example 2**  $k = 4, p = 3, s = 3$ . The distribution is shown in Figure 2(b). The transition table  $T$  is as follows. This demonstrates that the FSM for a problem can consist of disjoint cycles. Each processor sees at most one cycle of the FSM, depending on its start state.

STATE	NEXT	$\Delta\mathcal{M}$	$\Delta\mathcal{C}$	$\Delta\mathcal{O}$
0	3	3	0	3
1	1	4	1	0
2	2	4	1	0
3	0	1	1	-3



**Example 3**  $k = 4, p = 3, s = 5$ . The distribution is shown in Figure 2(c). The transition table  $T$  is as follows. The  $\Delta\mathcal{M}$  sequence has period  $k$ , and local memory accesses vary in stride.

STATE	NEXT	$\Delta\mathcal{M}$	$\Delta\mathcal{C}$	$\Delta\mathcal{O}$
0	3	7	1	3
1	2	9	2	1
2	0	2	1	-2
3	1	2	1	-2



### 2.1 The algorithm

Preparatory to finding a processor's  $\Delta\mathcal{M}$  sequence, we solve the problem of finding its starting memory location. Given a processor number  $m$ , regular section  $A(\ell : h : s)$ , and layout parameters  $p$  and  $k$ , we wish to find the first element (if any) of the regular section that resides on that processor. This is equivalent to finding the smallest nonnegative integer  $j$  such that  $\mathcal{P}(\ell + sj; p, k) = m$ .

Substituting the definition of  $\mathcal{P}$  from equation (1), we get

$$((\ell + sj) \bmod pk) \operatorname{div} k = m.$$

This reduces to

$$km \leq (\ell + sj) \bmod pk \leq k(m + 1) - 1$$

which is equivalent to finding an integer  $q$  such that

$$km - \ell \leq sj - pkq \leq km - \ell + k - 1. \quad (6)$$

Course	Processor 0				Processor 1				... Processor $p-1$			
0	0	...	$k-1$	$k$	...	$2k-1$	...	$(p-1)k$	...	$pk-1$		
1	$pk$	...	$(p+1)k-1$	$(p+1)k$	...	$(p+2)k-1$	...	$(2p-1)k$	...	$2pk-1$		
...	...	...	...	...	...	...	...	...	...	...		
Offset	0	...	$k-1$	0	...	$k-1$	...	0	...	$k-1$		

Local memory location	0	...	$k-1$	$k$	...	$2k-1$	...
Array index (processor 0)	0	...	$k-1$	$pk$	...	$(p+1)k-1$	...
Array index (processor 1)	$k$	...	$2k-1$	$(p+1)k$	...	$(p+2)k-1$	...
Array index (processor $p-1$ )	$(p-1)k$	...	$pk-1$	$(2p-1)k$	...	$2pk-1$	...

Figure 1: A one-level mapping. (a) Layout parameters. Processors run across the page, while memory is organized as a two-dimensional array of courses and offsets. Each box represents a memory cell, and the number in the cell indicates the array element it holds. (b) Layout of array elements in local processor memories. Each processor's memory is one row of the diagram.

Proc. 0	Proc. 1				Proc. 2				Proc. 3						
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31
32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47

Proc. 0	Proc. 1				Proc. 2						
0	1	2	3	4	5	6	7	8	9	10	11
12	13	14	15	16	17	18	19	20	21	22	23

Proc. 0	Proc. 1				Proc. 2						
0	1	2	3	4	5	6	7	8	9	10	11
12	13	14	15	16	17	18	19	20	21	22	23
24	25	26	27	28	29	30	31	32	33	34	35
36	37	38	39	40	41	42	43	44	45	46	47
48	49	50	51	52	53	54	55	56	57	58	59

Figure 2: One-level array mappings. The corresponding state tables are shown in Examples 1-3. (a)  $k = 4, p = 4, s = 3$ . (b)  $k = 4, p = 3, s = 3$ . (c)  $k = 4, p = 3, s = 5$ .

We rewrite inequality (6) as a set of  $k$  linear Diophantine equations in the variables  $j$  and  $q$ ,

$$\{sj - pkq = \lambda : km - \ell \leq \lambda \leq km - \ell + k - 1\}. \quad (7)$$

The equations in this set are to be solved independently (rather than simultaneously). Solutions exist for an individual equation if and only if  $\lambda$  is divisible by  $\text{GCD}(s, pk)$ .

The general solution to the linear Diophantine equation  $sj - pkq = \lambda$  is found as follows. Let  $d = \text{GCD}(s, pk) = \alpha s + \beta pk$ , with  $\alpha, \beta \in \mathbb{Z}$ . The quantities  $\alpha, \beta$ , and  $d$  are determined by the extended Euclid algorithm [6, page 325]. Then the general solution of the equation is the one-parameter family  $j = (\lambda\alpha + pk\gamma)/d$  and  $q = (-\lambda\beta + s\gamma)/d$ , for  $\gamma \in \mathbb{Z}$ .

For each solvable equation in set (7), we find the solution having the smallest nonnegative  $j$ . The minimum of these solutions gives us the starting index value  $\ell + js$ , which we then convert to the starting memory location using equation (4). We compute the ending memory location similarly, by finding the largest solutions no greater than  $h$ , and taking the maximum among these solutions.

Now we extend this idea to build the  $\Delta\mathcal{M}$  sequence for the processor. A sorted version of the set of smallest positive solutions determined above gives us the sequence of indices that will be

successively accessed by the processor. A linear scan through this sequence is sufficient to generate the  $\Delta\mathcal{M}$  sequence for the cycle of the FSM visible to processor  $m$ . If the FSM has multiple cycles, the local  $\Delta\mathcal{M}$  sequences may be different for different processors; otherwise, the sequences are cyclic shifts of one another. To avoid following the NEXT pointers at runtime, the  $\Delta\mathcal{M}$  table we compute contains the memory spacings in the order seen by the processor.

All of these elements are combined in Algorithm 1, which provides the solution for the memory address problem for a one-level mapping. The running time of Algorithm 1 is  $O(\log \min(s, pk)) + O(k \log k)$ , which reduces to  $O(\min(k \log k + \log s, k \log k + \log p))$ . As the output can be of size  $k$  in the worst case, any algorithm for this problem takes  $\Omega(k)$  time. Finding an algorithm with  $O(k)$  running time remains an open problem. Our implementation recognizes the following special cases that can be handled more quickly:  $p = 1$ ,  $s$  divides  $k$ ,  $k = 1$ ,  $pk$  divides  $s$ , a single course of blocks, and  $s$  divides  $pk$ .

Algorithm 1 can be simplified by noting that the right hand sides of successive solvable equations differ by  $d$ , and that the identity of the first solvable equation is easily determined from the value of  $(km - \ell) \bmod d$ . This removes the conditional from the loop in lines 5-12. Similar incremental computations can also be used to

simplify the floor and ceiling computations in that loop.  
The node code for each processor consists of the following loop.

```

base ← startmem; i ← 0
while (base ≤ lastmem) do
  compute with address base
  base ← base + ΔM[i]
  i ← (i + 1) mod length
enddo

```

**Algorithm 1 (Memory access sequence for a one-level mapping.)**

*Input:* Layout parameters  $(p, k)$ , regular section  $A(\ell : h : s)$ , processor number  $m$ .

*Output:* The  $\Delta\mathcal{M}$  table, the length of the table, starting memory location, ending memory location, and  $\sum \Delta\mathcal{M}$  for processor  $m$ .

*Method:*

```

1 integer ΔM[k], length
2 integer λ, i, start, end, d, α, β, lcourse, loffset, course,
  offset
3 integer loc[k+1], sorted[k+1]
4 (d, α, β) ← EXTENDED-EUCLID(s, pk); length ← 0;
  start ← h + 1; end ← ℓ - 1
5 for λ = km - ℓ to km - ℓ + k - 1 do
6   if λ mod d = 0 then /* Solvable equation */
7     loc[length] ← ℓ +  $\frac{\alpha}{d}(\lambda + pk \lceil \frac{-\lambda\alpha}{pk} \rceil)$ 
8     start ← min(start, loc[length])
9     end ← max(end, ℓ +  $\frac{\alpha}{d}(\lambda + pk \lfloor \frac{(h-1)d - \lambda\alpha}{pks} \rfloor))$ 
10    length ← length + 1
11  endif
12 enddo
13 if length = 0 or start > h then
14   return ΔM, ⊥, ⊥, ⊥, ⊥ /* No work for processor */
15 endif
16 sorted ← SORT(loc, length)
17 sorted[length] ← sorted[0] + pks/d; lcourse
  ← C(sorted[0]; p, k); loffset ← C(sorted[0]; p, k)
18 for i = 0 to length - 1 do
19   course ← C(sorted[i + 1]; p, k); offset
  ← C(sorted[i + 1]; p, k)
20   ΔM[i] ← k(course - lcourse) + (offset - loffset)
21   lcourse ← course; loffset ← offset
22 enddo
23 return ΔM, length, M(start; p, k), M(end; p, k),
  sk/GCD(s, pk)

```

We illustrate the actions of Algorithm 1 on Example 3 for processor 0. Here,  $p = 3, k = 4$ , and  $s = 5$ . The extended Euclid algorithm returns  $d = 1, \alpha = 5$ , and  $\beta = -2$ . Since each linear Diophantine equation is solvable, the loop in lines 5–12 is executed four times. At the end of this loop,  $loc = [0, 25, 50, 15]$ ,  $length = 4$ ,  $start = 0$ , and  $end = 50$ . After line 17,  $sorted = [0, 15, 25, 50, 60]$ . At the end of the loop in lines 18–22,  $\Delta\mathcal{M} = [7, 2, 9, 2]$ .

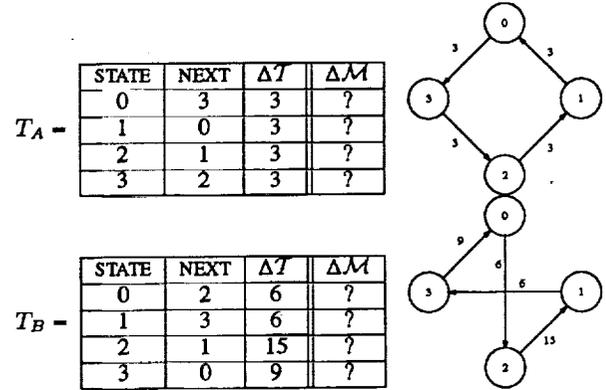
### 3 Two-level mappings

Two-level mappings are solved by two applications of Algorithm 1. As an example, let  $A$  be an array with element  $A(i)$  aligned with template cell  $3i$ , let the template be mapped with  $p = k = 4$ , and let the regular section be  $A(0 : 42 : 3)$ , which we abbreviate to  $B$ . The

picture is shown in Figure 3(a). The boxes represent template cells, the number  $j$  indicates the location of  $A(j)$ , and numbers in squares are part of the regular section  $A(0 : 42 : 3)$ . Two-level mappings encompass all mappings possible with the linguistic constructs of nested sections, affine alignment, and *cyclic*( $k$ ) distribution.

Unlike the one-level case, all template cells do not contain array elements. However, we allocate local storage only for the *occupied* cells of the template. Thus, the local memory storage for  $A$  in the example above is as shown in Figure 3(b). We handle this complication by first doing all our computations with respect to the template (the *local template space*), and then switching back to the local memory space in the final step. The picture in local template space is shown in Figure 3(c).

Note that  $A$  and  $B$  are both aligned to “regular sections” of the template:  $A(i)$  is aligned to template cell  $3i$ , and  $B(i)$  is aligned to template cell  $9i$ . We therefore create state tables for  $A$  and  $B$  with respect to the local template space using Algorithm 1. The “ $\Delta\mathcal{M}$ ” entries returned by the algorithm are actually spacings in local template space; we therefore rename this column  $\Delta\mathcal{T}$ . This gives the tables for  $A$  and  $B$  shown below.



Now all we need is to switch back to local memory space and fill in the  $\Delta\mathcal{M}$  entries.

As we traverse the state table  $T_A$  for any processor  $m$ , we access each element of  $A$  stored on that processor, in order. But since storage is allocated based on  $A$ , this translates to accessing consecutive locations in local memory. Thus the  $\Delta\mathcal{M}$  entries for  $T_A$  are all unity. Now it is a simple matter to find the  $\Delta\mathcal{M}$  entries for the desired section  $B$ . For instance, the first line of  $T_B$  says that if we are at an element of the section that has offset 0, the next element has offset 2 and is 6 template cells away. To find the corresponding  $\Delta\mathcal{M}$ , we start in state 0 of  $T_A$ , and walk the state table, accumulating  $\Delta\mathcal{T}$  and  $\Delta\mathcal{M}$  counts. We stop when we have accumulated the desired  $\Delta\mathcal{T}$  count. (It is easy to show that at this point we will have arrived at the correct offset as well.) The accumulated  $\Delta\mathcal{M}$  value is the desired entry for that line of  $T_B$ . This procedure results in the following table for  $B$ .<sup>1</sup>

STATE	NEXT	$\Delta\mathcal{T}$	$\Delta\mathcal{M}$
0	2	6	2
1	3	6	2
2	1	15	5
3	0	9	3

The algorithm as presented above requires  $O(k)$  time to compute each  $\Delta\mathcal{M}$  entry of  $T_B$ . Two optimizations reduce this complexity to  $O(1)$ , allowing the final fixup phase to run in  $O(k)$  time. First, we compute the sum of the  $\Delta\mathcal{T}$  and the  $\Delta\mathcal{M}$  around the state

<sup>1</sup>In this case,  $\Delta\mathcal{M}$  turns out to be  $\Delta\mathcal{T}$  divided by the alignment stride  $a = 3$ , but this is not true in general. Try the following case:  $a = 3, b = 0, p = 3, k = 8, s = 3$ .

(a)

Processor 0				Processor 1				Processor 2				Processor 3			
0			1			2			3			4			5
		6			7			8			9				10
	11			12			13			14			15		
16			17			18			19			20			21
		22			23			24			25			26	
	27			28			29			30			31		
32			33			34			35			36			37
		38			39			40			41			42	

(b)

Local memory location	0	1	2	3	4	5	6	7	8	9	10
Processor 0	0	1	6	11	16	17	22	27	32	33	38
Processor 1	2	7	12	13	18	23	28	29	34	39	
Processor 2	3	8	9	14	19	24	25	30	35	40	41
Processor 3	4	5	10	15	20	21	26	31	36	37	42

(c)

Local template cell	0	1	2	3	4	5	6	7	8	9	10	11	...
Processor 0	0			1			6			11			...
Processor 1			2			7			12			13	...
Processor 2		3			8			9			14		...
Processor 3	4			5			10			15			...

Figure 3: A two-level mapping of an array  $A$  aligned to the template with stride 3. (a) Layout of template cells. Numbers are array elements; boxed numbers are in the regular section  $A(0:42:3)$ . (b) Local memory storage for array. (c) Local template space for array.

Processor (0, 0)				Processor (0, 1)			
(0, 0)	(0, 1)	(0, 4)	(0, 5)	(0, 2)	(0, 3)	(0, 6)	(0, 7)
(1, 0)	(1, 1)	(1, 4)	(1, 5)	(1, 2)	(1, 3)	(1, 6)	(1, 7)
(2, 0)	(2, 1)	(2, 4)	(2, 5)	(2, 2)	(2, 3)	(2, 6)	(2, 7)
(9, 0)	(9, 1)	(9, 4)	(9, 5)	(9, 2)	(9, 3)	(9, 6)	(9, 7)
(10, 0)	(10, 1)	(10, 4)	(10, 5)	(10, 2)	(10, 3)	(10, 6)	(10, 7)
(11, 0)	(11, 1)	(11, 4)	(11, 5)	(11, 2)	(11, 3)	(11, 6)	(11, 7)
Processor (1, 0)				Processor (1, 1)			
(3, 0)	(3, 1)	(3, 4)	(3, 5)	(3, 2)	(3, 3)	(3, 6)	(3, 7)
(4, 0)	(4, 1)	(4, 4)	(4, 5)	(4, 2)	(4, 3)	(4, 6)	(4, 7)
(5, 0)	(5, 1)	(5, 4)	(5, 5)	(5, 2)	(5, 3)	(5, 6)	(5, 7)
(12, 0)	(12, 1)	(12, 4)	(12, 5)	(12, 2)	(12, 3)	(12, 6)	(12, 7)
(13, 0)	(13, 1)	(13, 4)	(13, 5)	(13, 2)	(13, 3)	(13, 6)	(13, 7)
(14, 0)	(14, 1)	(14, 4)	(14, 5)	(14, 2)	(14, 3)	(14, 6)	(14, 7)
Processor (2, 0)				Processor (2, 1)			
(6, 0)	(6, 1)	(6, 4)	(6, 5)	(6, 2)	(6, 3)	(6, 6)	(6, 7)
(7, 0)	(7, 1)	(7, 4)	(7, 5)	(7, 2)	(7, 3)	(7, 6)	(7, 7)
(8, 0)	(8, 1)	(8, 4)	(8, 5)	(8, 2)	(8, 3)	(8, 6)	(8, 7)
(15, 0)	(15, 1)	(15, 4)	(15, 5)	(15, 2)	(15, 3)	(15, 6)	(15, 7)
(16, 0)	(16, 1)	(16, 4)	(16, 5)	(16, 2)	(16, 3)	(16, 6)	(16, 7)
(17, 0)	(17, 1)	(17, 4)	(17, 5)	(17, 2)	(17, 3)	(17, 6)	(17, 7)

Figure 4: Mapping a two-dimensional array. Layout of array elements on processor arrangement with  $p = k = 3$  in the vertical dimension and  $p = k = 2$  in the horizontal dimension. The boxed pairs are  $A(0:17:2, 0:7:3)$ .

machine cycle, so that we can cast out any full cycles that must be traversed. Second, we tabulate the distance of each offset from a fixed origin (the *numarcs* table in Algorithm 2), allowing the computation of the distance between any two offset positions by two table lookups and some simple arithmetic. Consider the third line of  $T_B$ , which has a  $\Delta T$  of 15. Using the precomputed sums, we determine in constant time that this corresponds to one full cycle plus 3 more template cells, and that the full cycle contributes a memory offset of 4. Now all we need is to determine the number of memory cells encountered in going from offset position 2 to offset position 1 in  $T_A$ . We find this distance in constant time using the *numarcs* table. Algorithm 2 shows the method with these optimizations incorporated.

**Algorithm 2 (Memory access sequence for a two-level mapping.)**

*Input:* Alignment parameters  $(a, b)$ , layout parameters  $(p, k)$ , array length  $n$ , regular section  $A(l: h: s)$ , processor number  $m$ .

*Output:* The  $\Delta\mathcal{M}$  table, length of the table, starting memory location, ending memory location, and the length of  $A$  in processor  $m$ .

*Method:*

```

1 integer i, mm, numarcs[k]
2 integer  $\Delta\mathcal{M}[k]$ , length, startmem, lastmem
3 integer  $\Delta\mathcal{M}1[k]$ , length1, startmem1, lastmem1, memcycle1
4 integer  $\Delta\mathcal{M}2[k]$ , length2, startmem2, lastmem2, memcycle2

5 function residual(x) =
    length1 ·  $\lfloor (x - \text{startmem1}) / \text{memcycle1} \rfloor$ 
6 function major(x) = (numarcs[x mod k] -
    numarcs[startmem1 mod k]) mod length1
7 function realmem(x) = major(x) + residual(x)

8 ( $\Delta\mathcal{M}1$ , length1, startmem1, lastmem1, memcycle1)
    ← Algorithm 1(b, a(n - 1) + b, a, p, k, m)
9 ( $\Delta\mathcal{M}2$ , length2, startmem2, lastmem2, memcycle2)
    ← Algorithm 1(aℓ + b, ah + b, as, p, k, m)
10 if length1 = ⊥ or length2 = ⊥ then
11     return  $\Delta\mathcal{M}$ , ⊥, ⊥, ⊥, 0 /* No work for processor */
12 endif
13 mm ← startmem1
14 for i = 0 to length1 - 1 do
15     numarcs[mm mod k] ← i
16     mm ← mm +  $\Delta\mathcal{M}1[i]$ 
17 enddo
18 mm ← startmem2
19 for i = 0 to length2 - 1 do
20      $\Delta\mathcal{M}[i]$  ← realmem(mm +  $\Delta\mathcal{M}2[i]$ ) - realmem(mm)
21     mm ← mm +  $\Delta\mathcal{M}2[i]$ 
22 enddo
23 return  $\Delta\mathcal{M}$ , length2, realmem(startmem2),
    realmem(lastmem2), realmem(lastmem1) -
    realmem(startmem1) + 1

```

#### 4 Multidimensional arrays

So far, we have characterized the access sequence of a single processor by a start address and a memory stride pattern, the latter being generated by the transitions of a FSM. For a multidimensional array, each dimension will be characterized by a  $(p, k)$  pair and have a FSM associated with it, and each processor will execute

a set of nested loops. Instead of the start address being fixed as it was before, it too will be generated by the transitions of a FSM. For simplicity, we will deal with a two-dimensional array. The extension to higher dimensions is straightforward. An array dimension with  $p = 1$  is "mapped to memory".

Consider a one-level mapping of an array  $A(0: 17, 0: 7)$  with parameters  $(p_1, k_1) = (3, 3)$  in the first dimension and  $(p_2, k_2) = (2, 2)$  in the second dimension, and let the section of interest be  $A(0: 17: 2, 0: 7: 3)$ . The picture is shown in Figure 4.

Let each block be laid out in column-major order in local memory. Thus, the sequence of index values corresponding to successive memory elements of processor (0,0) is  $\{(0,0), (1,0), (2,0), (9,0), (10,0), (11,0), (0,1), \dots\}$ . Such a decomposition is used by Don-garra *et al.* for dense linear algebra codes on distributed-memory machines [3]. Then the tables corresponding to the two dimensions of  $A$  are

$$T_1 = \begin{array}{|c|c|c|} \hline \text{STATE} & \text{NEXT} & \Delta\mathcal{M} \\ \hline 0 & 2 & 2 \\ \hline 1 & 0 & 2 \\ \hline 2 & 1 & 2 \\ \hline \end{array}$$

and

$$T_2 = \begin{array}{|c|c|c|} \hline \text{STATE} & \text{NEXT} & \Delta\mathcal{M} \\ \hline 0 & 1 & 30 \\ \hline 1 & 0 & 6 \\ \hline \end{array}$$

(Had the section been larger in the second dimension, the column accessed by processor (0,0) immediately after column 0 would be column 9. The elements (0,0) and (0,9) would be 30 locations apart in local memory. This explains the  $\Delta\mathcal{M}$  entry in the first row of  $T_2$ .)

The FSM  $T_2$  is used to generate memory spacings for successive base addresses, while the FSM  $T_1$  is used to generate memory spacings between successive elements with the same base address. Algorithm 3 computes these FSMs.

**Algorithm 3 (Memory access sequence for a multidimensional array with a two-level mapping.)**

*Input:* Number of dimensions  $d$ ; alignment parameters  $(a_j, b_j)$ , layout parameters  $(p_j, k_j)$ , dimension length  $n_j$ , regular section  $A(l_j: h_j: s_j)$ , processor number  $m_j$ , for each dimension  $j$  ( $1 \leq j \leq d$ ).

*Output:*  $\Delta\mathcal{M}_j$ , length $_j$ , startmem $_j$ , lastmem $_j$ , ( $1 \leq j \leq d$ ).

*Method:*

```

1 integer j, mem, m, k
2 mem ← 1
3 for j = 1 to d do
4     ( $\Delta\mathcal{M}_j$ , length $_j$ , startmem $_j$ , lastmem $_j$ , m)
        ← Algorithm 2( $a_j, b_j, p_j, k_j, n_j, l_j, h_j, s_j, m_j$ )
5     startmem $_j$  ← startmem $_j$  · mem
6     lastmem $_j$  ← lastmem $_j$  · mem
7     for k = 0 to length $_j$  - 1 do
8          $\Delta\mathcal{M}_j[k]$  ←  $\Delta\mathcal{M}_j[k]$  · mem
9     enddo
10    mem ← mem · m
11    enddo
12    return  $\Delta\mathcal{M}_j$ , length $_j$ , startmem $_j$ , lastmem $_j$ ;
        ( $1 \leq j \leq d$ )

```

For the two-dimensional case, each processor executes the following doubly-nested loop. Higher-dimensional cases follow the same pattern, with one loop for each dimension.

```

base2 ← startmem2
i2 ← 0
while base2 ≤ lastmem2 do
  base1 ← startmem1
  i1 ← 0
  while base1 ≤ lastmem1 do
    compute with address base2 + base1
    base1 ← base1 + ΔM1[i1]
    i1 ← (i1 + 1) mod length1
  enddo
  base2 ← base2 + ΔM2[i2]
  i2 ← (i2 + 1) mod length2
enddo

```

## 5 Generating communication sets

The set of messages that a given processor  $m$  must send while performing its share of the computation is called its *communication set*. We now extend the FSM approach to generate communication sets for regular communication patterns such as shifts and stride changes. In this section, we assume a one-dimensional array with a single-level mapping. Extensions to multilevel and multidimensional cases are as in the address generation problem.

We adopt Koelbel's execution model for parallel loops [7], in which a processor goes through four steps:

1. Send those data items it holds that are required by other processors.
2. Perform local iterations, *i.e.*, iterations for which it holds all data.
3. Receive data items from other processors.
4. Perform nonlocal iterations.

We combine the final two steps into a receive-execute loop.

Our approach puts all the intelligence on the sending side, keeping the receive-execute loop simple. We wish to minimize the number of messages sent and maximize cache benefits. To this end, each processor makes a single pass over the data in its local memory using the FSM technique described above, figuring out the destination of each data element and building messages. The messages contain (address, value) pairs which the receiving processor uses to perform its computations. All the data that must go from processor  $i$  to processor  $j$  are communicated in a single message. The data items corresponding to the local iterations are conceptually sent by the processor to itself. Of course, this message is never sent; rather, the local iterations are performed out of the message buffer. The idea is similar to the inspector-executor model [10] for irregular loops.

### 5.1 Shifts

Consider the computation

$$A(0:h:s) = A(0:h:s) + A(\ell:\ell+h:s),$$

where the regular section  $A(\ell:\ell+h:s)$  is to be moved to align with  $A(0:h:s)$ . (For simplicity, assume  $\ell > 0$ .) In this case, a processor communicates with at most two processors. We show how to augment the FSM with a  $\Delta\mathcal{P}$  and a  $\Delta\mathcal{L}$  column, such that adding these quantities to the processor and memory location of the current element gives the remote processor number and the memory location in the remote processor of the element with which it interacts.

Suppose we know that element  $A(i)$  is located on processor  $m$  at memory location  $\mathcal{M}(i; p, k)$  with offset  $O_i = \mathcal{O}(i; p, k)$ . Then the problem is to find the processor and local memory location of the element  $A(i - \ell)$ . Let  $\ell = q \cdot pk + r$ , and  $\Delta\mathcal{P} = \lceil (r - O_i)/k \rceil$ . Then  $\mathcal{P}(i - \ell; p, k) = (m - \Delta\mathcal{P} + p) \bmod p$ . Since  $0 \leq O_i < k$ ,  $\Delta\mathcal{P}$  can assume only two values as the offset  $O_i$  varies.

We also define the quantity  $\Delta\mathcal{L}$  such that  $\mathcal{M}(i - \ell; p, k) = \mathcal{M}(i; p, k) + \Delta\mathcal{L}(O_i)$ .

$$\Delta\mathcal{L}(O_i) = ((O_i - r) \bmod k) - O_i - kq - \eta,$$

$$\eta = \begin{cases} k & \text{if } (mk - r + O_i) < 0, \\ 0 & \text{otherwise.} \end{cases}$$

### 5.2 Stride changes

Now consider the communication required to move the regular section  $A(0:s_2n:s_2)$  to align with the regular section  $A(0:s_1n:s_1)$ . This is harder than a shift; there is no simple lookup technique for generating the communication sets. The problem is that the pattern of destination processors can have period longer than  $k$ . The only fact that appears to be true in this case is that the sum of the periods over all processors divides  $pk$ . We therefore resort to a simple technique.

We generate the FSM to reflect the storage properties of the  $s_2$ -section. Now, the  $j^{\text{th}}$  element of this section interacts with index  $i = j \cdot s_1$  of  $A$ . Let  $q_1 = i \text{ div } pk$ ,  $r_1 = i \bmod pk$ ,  $q_2 = r_1 \text{ div } k$ , and  $r_2 = r_1 \bmod k$ . Then the desired element is stored in memory location  $q_1 \cdot k + r_2$  on processor  $q_2$ .

Many modern RISC microprocessors perform integer division and remainder operations in software, which are rather slow. The equations above can be rewritten to require only two integer divisions, which substantially speeds up the implementation. If either  $p$  or  $k$  is a power of two, the computation can be sped up further by using shifts and bit masks instead of the division operations. Our implementation incorporates these optimizations.

### 5.3 Termination

Each processor must determine when it has completed all its operations. As we do not send messages between every pair of processors, we cannot use the number of messages to determine termination. However, it is easy to figure out the total number of operations that processor  $m$  will execute (by a variant of Algorithm 1). A simple termination test initializes a counter to this value, decrements it by the number of array elements received in each message, and exits when the counter reaches zero.

## 6 Experimental results

In this section we present experimental results of implementing the algorithms above. We measure both the preprocessing cost of table generation and the runtime loop overhead.

Our computational kernel is the DAXPY operation

$$y(1:n:s) = y(1:n:s) + a * x(1:n:s).$$

The ratio of memory operations to computation is high in this kernel, making it unfavorable for our algorithms. The overhead due to irregular memory access patterns would decrease with more computation in the loop. All experiments reported in this section were done on a Sparcstation 2, for which Dongarra [2] reports a  $100 \times 100$  double-precision LINPACK performance of 4.0 Mflops. (We measured 3.3 Mflops.) (A repetition of the experiments on a single processor of an iPSC/860 showed similar characteristics.) We used the GNU C compiler, with optimization at the -O2 level.

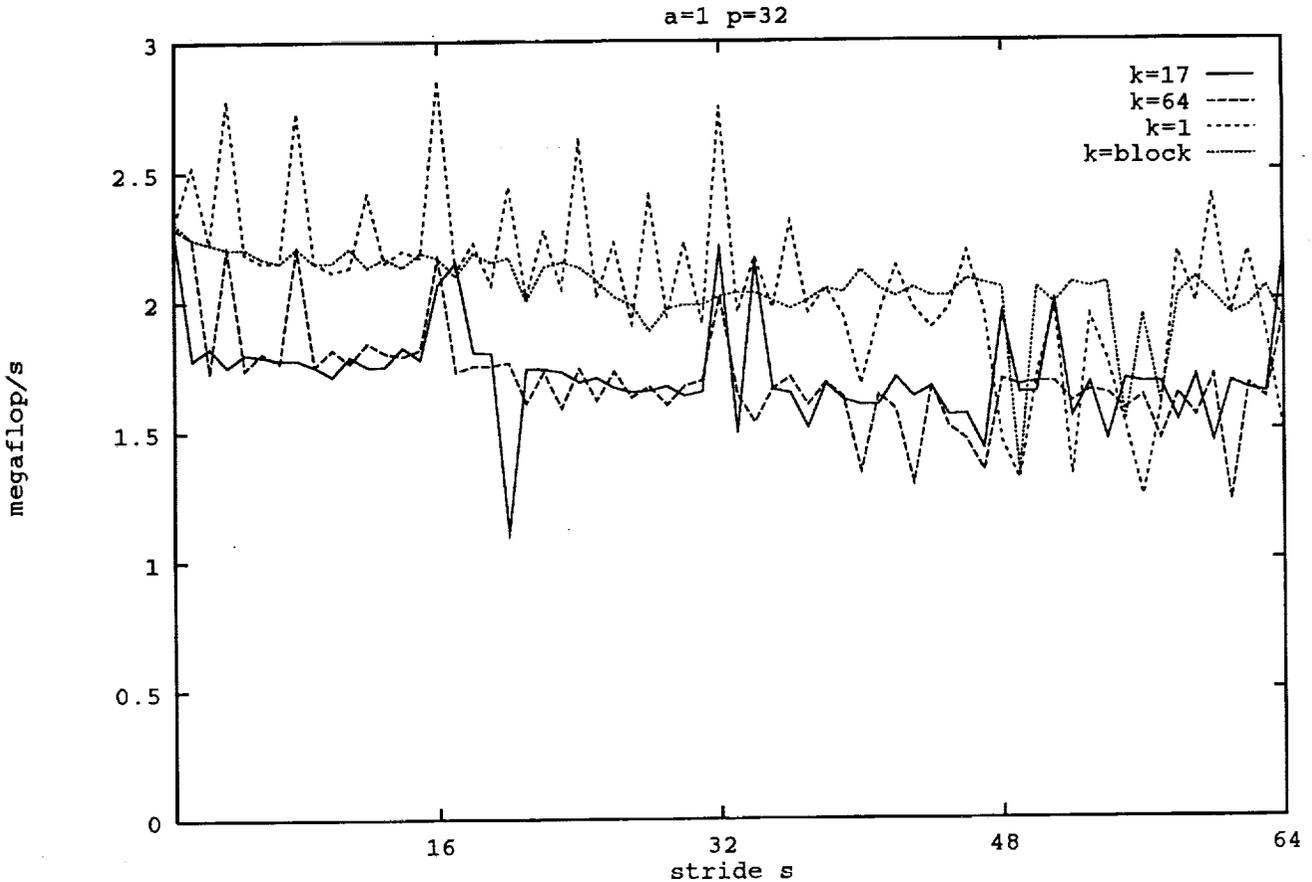


Figure 5: Performance of DAXPY for various block sizes  $k$ .

A higher optimization level increased compilation time without improving the quality of the generated code. With  $s = 1$  the loop compiled into 9 instructions, while the general case with table lookup compiled into 17 instructions. Our measured performance for the  $s = 1$  case was 2.0 Mflops for long vectors.

### 6.1 Local address generation

We have seven parameters for a one-dimensional array: two alignment parameters  $a$  and  $b$ , two layout parameters  $p$  and  $k$ , and three section parameters  $\ell$ ,  $h$ , and  $s$ . This is too large a space over which to present results. We now explain how we reduce the number of variables.

- **Alignment stride  $a$ :** For any constant  $c$ , the memory layout for parameters  $ca$  and  $ck$  is exactly the same as that for  $a$  and  $k$ . Thus we experiment only with  $a = 1$ ; the effect of changing  $a$  can be simulated by changing  $k$  instead.
- **Alignment offset  $b$ :** This is irrelevant for large sections. For small sections where end effects could be important, we observed that changes in  $b$  affected results by less than 5%. We therefore take  $b = 0$ .
- **Section lower bound  $\ell$ :** Again, this is irrelevant for large sections and has negligible effect even for small sections. We therefore take  $\ell = 0$ .

- **Section upper bound  $h$ :** This affects the number of elements in the DAXPY, which in turn affects the tradeoff between loop startup and iteration time in a single processor. We vary  $h$  along with  $s$ , keeping  $h/s$  fixed, so that all our experiments do the same amount of work regardless of stride.
- **Section stride  $s$ :** This is the  $x$ -axis of our plots, as one of the two independent variables.
- **Number of processors  $p$ :** We run our experiments with a fixed number of processors. Varying the number of processors has the effect of scaling the plots.
- **Block size  $k$ :** This varies from plot to plot, as the second of the two independent variables.

We use the execution time per element of the DAXPY loop as our measure of performance. For a given set of parameters, we compute the tables for each processor and time the execution of the distributed loop. Let  $\tau_m$  be the time taken by processor  $m$  to execute the  $n_m$  iterations it receives. Then the execution time per element is

$$t_e = \max_m \frac{\tau_m}{n_m}.$$

This factors out the effect of load imbalance.

We plot  $2/t_e$  (the megaflops per second per processor) against stride  $s$  for various block sizes  $k$ . The plots are shown in Figure 5. The curves are the sum of two effects.

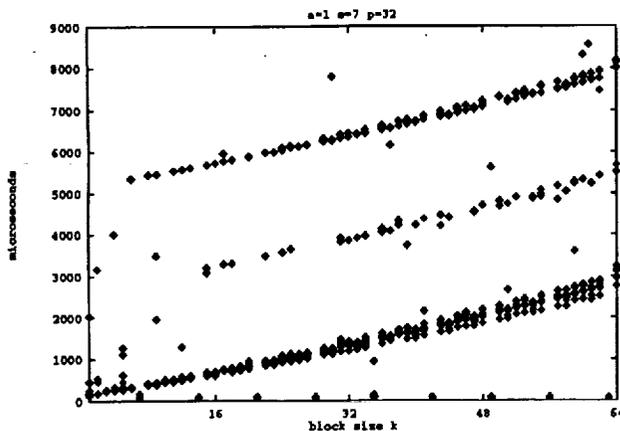


Figure 6: Preprocessing overhead for memory stride computations. The graph corresponds to an alignment stride of 1. The multiple points for each block size show the processing times for each of 32 processors.

1. **Baseline:** The lower envelope is a smooth, relatively flat function of  $s$ . The value of the envelope is approximately the asymptotic performance for the single-processor DAXPY.
2. **Jiggles:** The plot jiggles up and down some, presumably due to cache stride effects [1].

Although the general loop has almost twice as many instructions as the unit-stride loop, the difference in performance is much less. This is because the execution time is dominated by the floating-point operations and references to operands in memory. The references to the  $\Delta\mathcal{M}$  array usually hit in cache, and the additional loop control instructions probably execute in overlapped fashion.

The preprocessing time required to compute the  $\Delta\mathcal{M}$  tables is shown in Figure 6, where we plot the running time of Algorithm 1 against block size  $k$ . The running time is approximately  $O(k)$ , with certain special cases being handled much faster. Experimental observation confirmed that the running times for Algorithm 2 were no more than twice those of Algorithm 1.

If the  $\Delta\mathcal{M}$  sequence is known at compile time, unrolling the loop bodies avoids the indirection into the  $\Delta\mathcal{M}$  array. Figure 7 shows that the benefits of this are mixed. While loop unrolling removes the indirection, it also increases the size of the loop body; thus the benefits of a simpler loop body may be undone by misses in the instruction cache caused by the larger loop body. This is a factor for larger block sizes, since the amount of unrolling may be as large as  $k$ . We recommend unrolling only if the  $\Delta\mathcal{M}$  sequences are short, or can be compacted, *e.g.*, by run-length encoding.

## 6.2 Communication set generation

For the communication set generation algorithms, we measured the time required to generate the FSM tables and the time required to marshal the array elements into messages. We did not measure the actual communication time, since that varies widely among machines. The results are presented in Table 2.

## 7 Conclusions

Our table-driven approach to the generation of local memory addresses and communication sets unifies and subsumes previous solutions to these problems. The tables required in our method are

$k$	$s$	$\ell$	FSM gen. time ( $\mu\text{s}$ )	Run time/elt ( $\mu\text{s}$ )
1	3	10	360	3.3
17	3	10	1831	3.1
17	3	40	1870	3.1
64	3	10	6632	3.2
64	3	100	6508	3.2
blk	3	10	408	3.2

$k$	$s_1$	$s_2$	FSM gen. time ( $\mu\text{s}$ )	Run time/elt ( $\mu\text{s}$ )
1	3	5	360	5.3
17	3	5	1891	14.3
64	3	5	6542	5.4
blk	3	5	400	11.3

Table 2: Performance of communication set generation algorithms. All numbers are for 32 processors. (a) Shift communication: the shift distance is  $\ell$ . (b) Change of stride: the stride is changed from  $s_2$  to  $s_1$ .

easy to compute, and the runtime overheads due to indirect addressing of data are small. Our data access patterns exploit temporal locality and make effective use of cache. We also support blocking of messages to reduce message startup overheads.

Our prototype implementation has demonstrated our performance claims. An optimized implementation of our techniques should deliver close to maximum performance for general *cyclic*( $k$ ) distributions of arrays.

Our algorithms are fully parallel in that each processor of a parallel machine can generate its tables independently. Finding a fully parallel algorithm running in  $O(k)$  parallel time or with  $O(p+k)$  total work remains an open problem.

## Acknowledgments

J. Ramanujam proofread an early draft and suggested the simplifications to Algorithm 1. We also thank the referees for their helpful comments.

## References

- [1] BAILEY, D. H. Unfavorable strides in cache memory systems. RNR Technical Report RNR-92-015, Numerical Aerodynamic Simulation Systems Division, NASA Ames Research Center, Moffett Field, CA, May 1992.
- [2] DONGARRA, J. J. Performance of various computers using standard linear equations software. *Computer Architecture News* 20, 3 (June 1992), 22–44.
- [3] DONGARRA, J. J., VAN DE GEIJN, R., AND WALKER, D. W. A look at scalable dense linear algebra libraries. In *Proceedings of the Scalable High Performance Computing Conference* (Williamsburg, VA, Apr. 1992), IEEE Computer Society Press, pp. 372–379. Also available as technical report ORNL/TM-12126 from Oak Ridge National Laboratory.
- [4] FOX, G. C., HIRANANDANI, S., KENNEDY, K., KOELBEL, C., KREMER, U., TSENG, C.-W., AND WU, M.-Y. Fortran D language specification. Tech. Rep. Rice COMP TR90-141, Department of Computer Science, Rice University, Houston, TX, Dec. 1990.

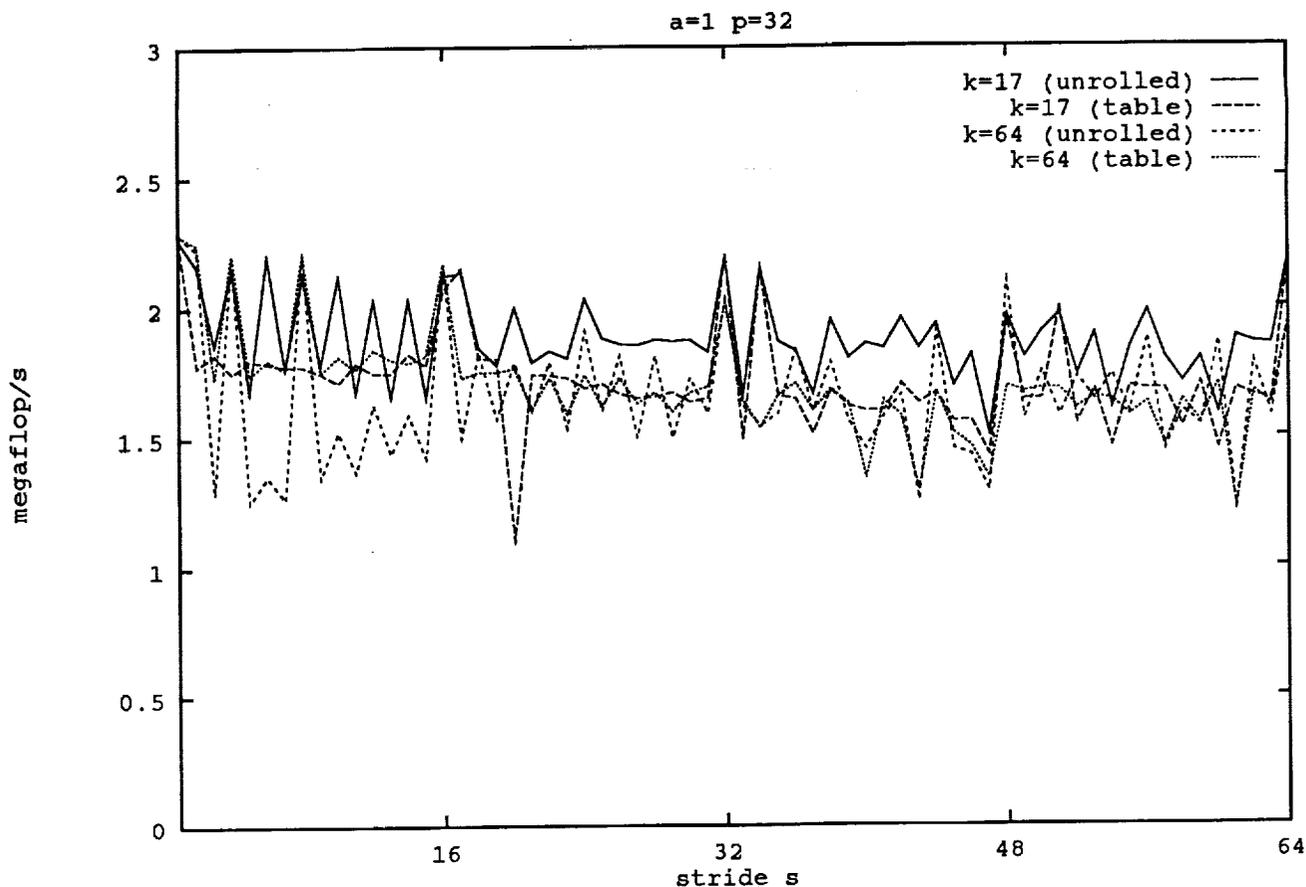


Figure 7: Effect of loop unrolling on performance of DAXPY for various block sizes  $k$ .

- [5] HIGH PERFORMANCE FORTRAN FORUM. High Performance Fortran language specification version 0.4. Draft, Nov. 1992. Also available as technical report CRPC-TR 92225, Center for Research on Parallel Computation, Rice University.
- [6] KNUTH, D. E. *Seminumerical Algorithms*, second ed., vol. 2 of *The Art of Computer Programming*. Addison-Wesley Publishing Company, Reading, MA, 1981.
- [7] KOELBEL, C. Compile-time generation of regular communication patterns. In *Proceedings of Supercomputing'91* (Albuquerque, NM, Nov. 1991), pp. 101-110.
- [8] KOELBEL, C., AND MEHROTRA, P. Compiling global namespace parallel loops for distributed execution. *IEEE Transactions on Parallel and Distributed Systems* 2, 4 (Oct. 1991), 440-451.
- [9] MACDONALD, T., PASE, D., AND MELTZER, A. Addressing in Cray Research's MPP Fortran. In *Proceedings of the Third Workshop on Compilers for Parallel Computers* (Vienna, Austria, July 1992), Austrian Center for Parallel Computation, pp. 161-172.
- [10] MIRCHANDANEY, R., SALTZ, J. H., SMITH, R. M., CROWLEY, K., AND NICOL, D. M. Principles of runtime support for parallel processors. In *Third International Conference on Supercomputing* (July 1988), ACM Press, pp. 140-152.