



HAL
open science

How to optimize residual communications ?

Michèle Dion, Cyril Randriamaro, Yves Robert

► **To cite this version:**

Michèle Dion, Cyril Randriamaro, Yves Robert. How to optimize residual communications?. [Research Report] LIP RR-1995-27, Laboratoire de l'informatique du parallélisme. 1995, 2+26p. hal-02101780

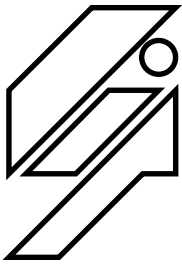
HAL Id: hal-02101780

<https://hal-lara.archives-ouvertes.fr/hal-02101780v1>

Submitted on 17 Apr 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Laboratoire de l'Informatique du Parallélisme

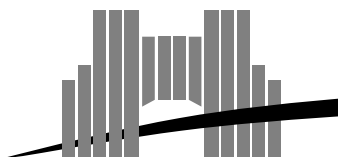
Ecole Normale Supérieure de Lyon
Unité de recherche associée au CNRS n°1398

How to optimize residual communications ?

Michèle Dion
Cyril Randriamaro
Yves Robert

October 1995

Research Report N° 95-27



Ecole Normale Supérieure de Lyon

46 Allée d'Italie, 69364 Lyon Cedex 07, France

Téléphone : (+33) 72.72.80.00 Télécopieur : (+33) 72.72.80.80

Adresse électronique : lip@lip.ens-lyon.fr

How to optimize residual communications ?

Michèle Dion
Cyril Randriamaro
Yves Robert

October 1995

Abstract

Minimizing communications when mapping affine loop nests onto distributed memory parallel computers has already drawn a lot of attention. This paper focuses on the next step: as it is generally impossible to obtain a communication-free (or local) mapping, how to optimize the residual communications? We explain how to take advantage of macro-communications such as broadcasts, scatters, gathers or reductions or how to decompose general affine communications into simpler ones that can be performed more efficiently. We finally give a two-step heuristic that summarizes our approach: first minimize the number of nonlocal communications, then optimize residual affine communications using macro-communications or decompositions.

Keywords: mapping, affine loop nests, distributed memory parallel computers, minimizing communications, macro-communications

Résumé

Minimiser les communications lors du placement d'un nid de boucles affine sur une machine parallèle à mémoire distribuée a déjà été l'objet de beaucoup d'attention. Dans ce rapport, nous nous intéressons à l'étape suivante : puisqu'il est en général impossible d'obtenir un placement sans communication, comment optimiser les communications résiduelles ? Nous expliquons comment tirer avantage de macro-communications (diffusions, distributions, rassemblements ou réductions) ou comment décomposer des communications affines en communications plus simples et plus efficaces. Nous donnons finalement une heuristique en 2 étapes qui résume notre approche : d'abord minimiser le nombre de communications non locales , puis optimiser les communications résiduelles en utilisant des macro-communications ou des décompositions.

Mots-clés: placement, nids de boucles affines, machines parallèles à mémoire distribuée, minimisation des communications, macro-communications

How to optimize residual communications ?

Michèle Dion, Cyril Randriamaro and Yves Robert *

Laboratoire LIP, URA CNRS 1398

Ecole Normale Supérieure de Lyon, F - 69364 LYON Cedex 07

e-mail: [Michele.Dion,Cyril.Randriamaro,Yves.Robert]@lip.ens-lyon.fr

Abstract

Minimizing communications when mapping affine loop nests onto distributed memory parallel computers has already drawn a lot of attention. This paper focuses on the next step: as it is generally impossible to obtain a communication-free (or local) mapping, how to optimize the residual communications? We explain how to take advantage of macro-communications such as broadcasts, scatters, gathers or reductions or how to decompose general affine communications into simpler ones that can be performed more efficiently. We finally give a two-step heuristic that summarizes our approach: first minimize the number of nonlocal communications, then optimize residual affine communications using macro-communications or decompositions.

1 Introduction

This paper deals with the problem of mapping affine loop nests onto Distributed Memory Parallel Computers (DMPCs). Because communication is very expensive in DMPCs, how to distribute data arrays and computations to processors is a key factor to performance.

The computations described in this paper are general non-perfect loop nests (or multiple loop nests) with uniform or affine dependences. Mapping such loop nests onto virtual processors with the aim of “minimizing” in some sense the communication volume or number is known as the alignment problem, which has drawn a lot of attention [15, 9, 13, 1, 4, 20, 12, 3, 6].

More precisely, given a loop nest, the goal is to assign a virtual processor, i.e. a location on a virtual processor grid, to each array element and to each computation. Arrays and computations are thus aligned and projected onto the virtual processor grid. This grid is then folded onto a physical grid, most often of much smaller size in each dimension. Languages like *HPF* provide *BLOCK* or *CYCLIC* distributions for such a folding.

There is a natural objective function for the mapping problem which has been extensively used in the literature: the aim is to minimize the number of non-local communications that will remain after a suitable alignment has been found. Some authors even look for a *communication-free* mapping for which there remains *NO* nonlocal communication. In fact, it is always possible to achieve a communication-free mapping ... provided we are prepared to use a virtual grid of dimension 0, i.e. a single processor ! Obviously, the larger the dimension of the target virtual grid, the larger the number of residual nonlocal communications.

*Supported by the ESPRIT Basic Research Action 6632 *NANA2* of the European Economic Community and by the CNRS-INRIA Project *ReMaP*.

Experimental evidence shows that communication-free mappings are very unlikely to be achieved. Think of elementary kernels as simple as a matrix-matrix product or a Gaussian elimination procedure, there is no way to map such kernels onto 2D- or even 1D-grids without residual communications. A natural question arises: is there a way to “optimize” in some sense the communications that remain ?

Platonoff [19] gives a strong motivation to answer the question. Experimenting on a $CM - 5$ with 32 processors, he compared various communication times. He observed the ratios reported in Table 1. Table 1 clearly shows that the $CM - 5$ has facilities for some usual macro-communications such as a broadcast or a reduction. Also, translations are much more efficient than general (affine) communications.

Reduction	Broadcast	Translation	General communication
1	1.5	2.5	90

Table 1: Comparing execution times for different data movements on a CM-5

Experimenting with the Intel Paragon gives interesting results too: communication conflicts are generated by serial messages on a single link at the same time. Thus a general communication cannot be efficiently executed by simply letting all processors send their messages simultaneously. Decomposing a general communication into a small sequence of communication parallel to the grid axes (horizontal and vertical) proves much faster, as illustrated in Section 4.1.

Being prepared to have some residual communications, we adopt the following strategy:

1. Zero out as many non-local communications as possible. To this purpose, we use an heuristic modified from Dion and Robert [6]. This heuristic is based upon the *access graph* and will be explained in Section 2. We weight the edges of the access graph to take communication volume into account.
2. For remaining communications, we explore the following two possibilities (both can be implemented simultaneously):
 - (a) try to find a mapping such that (at least) one of the residual communications is a macro-communication (in addition to reduction and broadcast, we could have a scatter or a gather, for instance).
 - (b) try to find a mapping such that (at least) one of the residual communications can be decomposed into more simple and efficient data movements (such as horizontal and vertical communications).

The paper is organized as follows: in Section 2, we introduce a motivating example that we will work out throughout the paper. We informally explain our modified heuristic to zero out as many communications as possible.

We explain how to take advantage of macro-communications in Section 3. We explain how to decompose general affine communications into simpler ones in Section 4. Altogether, our complete heuristic for the mapping problem is summarized in Section 5. Section 6 is devoted to a survey of related work in the literature, together with some discussion and comparisons. Finally, concluding remarks are stated in Section 7.

2 Motivating example

In this section, we informally explain our new approach on an example. First we apply an heuristic modified from [6] to minimize the number of nonlocal communications. Then, we try to make the residual nonlocal communications “efficient”.

2.1 Example

Consider the following non-perfect affine loop nest, where a is a $2D$ -array, and b and c are $3D$ -arrays:

Example 1

```

for  $i = 0$  to  $N$  do
  for  $j = 0$  to  $M$  do
    { Statement  $S_1$  }  $b(F_1 \cdot (i, j)^t + c_1) = g_1(a(F_2 \cdot (i, j)^t + c_2), a(F_3 \cdot (i, j)^t + c_3),$ 
                                      $c(F_4 \cdot (i, j)^t + c_4))$ 
    for  $k = 0$  to  $N + M$ 
      { Statement  $S_2$  }  $b(F_5 \cdot (i, j, k)^t + c_5) = g_2(a(F_6 \cdot (i, j, k)^t + c_6))$ 
      { Statement  $S_3$  }  $c(F_7 \cdot (i, j, k)^t + c_7) = g_3(a(F_8 \cdot (i, j, k)^t + c_8))$ 
    endfor
  endfor
endfor

```

where

$$\begin{aligned}
 F_1 &= \begin{pmatrix} 1 & 1 \\ 0 & 0 \\ 0 & -1 \end{pmatrix} \text{ and } c_1 = \begin{pmatrix} -2 \\ 1 \\ 3 \end{pmatrix}, F_2 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \text{ and } c_2 = \begin{pmatrix} -1 \\ 6 \end{pmatrix}, F_3 = \begin{pmatrix} 1 & 1 \\ -4 & -5 \end{pmatrix} \\
 \text{and } c_3 &= \begin{pmatrix} 0 \\ 0 \end{pmatrix}, F_4 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & -1 \end{pmatrix} \text{ and } c_4 = \begin{pmatrix} -1 \\ 3 \\ 0 \end{pmatrix}, F_5 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = Id_3 \text{ and } c_5 = \\
 \begin{pmatrix} 2 \\ 5 \\ -3 \end{pmatrix}, F_6 = \begin{pmatrix} 0 & 1 & -1 \\ 1 & -1 & 1 \end{pmatrix} \text{ and } c_6 = \begin{pmatrix} 1 \\ -2 \end{pmatrix}, F_7 = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & -1 \end{pmatrix} \text{ and } c_7 = \begin{pmatrix} 0 \\ -1 \\ 0 \end{pmatrix}, \\
 F_8 = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \end{pmatrix} \text{ and } c_8 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.
 \end{aligned}$$

Here, g_1 , g_2 and g_3 are arbitrary functions. The loop nest is an *affine loop nest* because all array references are affine functions of the loop indices. There are no data dependences in the nest (check this with Tiny [22] for instance) can be used, for example, to perform dependence analysis), all loops are DOALL loops, hence all computations can be executed at the same time-step.

Mapping the loop nest onto a m -dimensional virtual processor space consists in determining allocation matrices for each statement and for each array. *Affine allocation functions* are defined as in [1, 6]:

1. for each statement S of depth d : $\text{alloc}_S(I) = M_S I + \alpha_S$, where I is the iteration vector ($I = (i, j)^t$ for statement S_1 and $I = (i, j, k)^t$ for statements S_2 and S_3), M_S is a $m \times d$ matrix and α_S is a m -vector.
2. for each array x of dimension q_x : $\text{alloc}_x(I) = M_x I + \alpha_x$, where M_x is a $m \times q_x$ matrix and α_x is a m -vector.

By this mapping, statement instance $S(I)$ is assigned to virtual processor $M_S I + \alpha_S$ and array element $x(I)$ is stored in the local memory of processor $M_x I + \alpha_x$. For instance, in statement S_1 , the value $a(F_2(i, j)^t + c_2)$ has to be read in the memory of the processor

$$\text{alloc}_a(F_2(i, j)^t + c_2) = M_a(F_2(i, j)^t + c_2)I + \alpha_a$$

and sent to the processor in charge of the computation $S_1((i, j)^t)$, namely processor

$$\text{alloc}_{S_1}((i, j)^t) = M_{S_1}(i, j)^t + \alpha_{S_1}.$$

This results in a communication of length δ_{a, S_1} equal to the distance between $a(F_2(i, j)^t + c_2)$ and $S_1(i, j)$. We have

$$\begin{aligned} \delta_{a, S_1} &= \text{alloc}_{S_1}((i, j)^t) - \text{alloc}_a(F_2(i, j)^t + c_2) \\ \delta_{a, S_1} &= M_{S_1}(i, j)^t + \alpha_{S_1} - (M_a(F_2(i, j)^t + c_2) + \alpha_a) \\ \delta_{a, S_1} &= (M_{S_1} - M_a F_2)(i, j)^t - M_a c_2 + \alpha_{S_1} - \alpha_a. \end{aligned}$$

The expression δ_{a, S_1} contains a *nonlocal* term $(M_{S_1} - M_a F_2)(i, j)^t$ which depends upon the iteration vector $(i, j)^t$, and a *local* term $\alpha_{S_1} - M_a c_2 - \alpha_a$. The *nonlocal* term corresponds to irregular patterns of communication, whose size can grow over the whole processor space. It is clearly the main factor affecting performance. On the other hand, the *local* term corresponds to regular fixed-size communications that can be performed efficiently onto DMPCs.

Zeroing out the nonlocal term is of course the main goal of the mapping optimization process, as recognized by [1] and [20]. A *communication-local* mapping is a mapping where all nonlocal terms have been zeroed out. However, it is generally impossible to obtain a *communication-local* mapping with the required dimension for the target virtual architecture. Then, another objective of the mapping process is to make efficient those nonlocal communications that cannot be zeroed out. We will try for example to derive macro-communications such as broadcasts, reductions, scatters, ..., or to decompose general communications into communications parallel to the axes of the target processor space.

Back to our example, in statement S_1 and for array a , we see that zeroing out the nonlocal term amounts to choose allocation matrices M_{S_1} and M_a so that the equation $M_{S_1} - M_a F_2 = 0$ is satisfied. In the same way, the value $b(i + j - 2, 1, 3 - j)$ has to be written after the computation $S_1((i, j)^t)$. This results in a communication:

$$\begin{aligned} \delta_{S_1, b} &= \text{alloc}_b(F_1(i, j)^t + c_1) - \text{alloc}_{S_1}((i, j)^t) \\ \delta_{S_1, b} &= (M_b F_1 - M_{S_1})(i, j)^t + M_b c_1 + \alpha_b - \alpha_{S_1} \end{aligned}$$

Again, to zero out the nonlocal term, we have to fulfill equation $M_b F_1 - M_{S_1} = 0$.

More generally, if an array x is read or written in a statement S with an access matrix F , to zero out the nonlocal part of the communication due to this access, we must satisfy the matrix equation

$$M_S = M_x F.$$

There are 8 such equations in our example.

2.2 Zeroing out nonlocal communications

The primary goal is to zero out as many nonlocal terms as possible. That is to say: given M_S (resp. M_x) of full rank¹, we want to find M_x (resp. M_S) of full rank such as $M_S = M_x F$. How to solve such an equation is explained below.

¹We search for full rank allocation matrices: otherwise the target m -dimensional processor space would not be fully utilized.

2.2.1 Solving $M_S = M_x F$

Consider the nonlocal term linking statement S of depth d and array x of dimension q_x : the equation is $M_S = M_x F$, where M_S is a $m \times d$ matrix, and M_x a $m \times q_x$ matrix. The matrix F is of dimension $q_x \times d$ and we consider only full rank access matrices. We target a m -dimensional processor space and we assume that $m \leq d$ and $m \leq q_x$ (we have chosen to deal only with the communications such that $m \leq d$ and $m \leq q_x$ as they represent the core of the computations and data elements to be distributed). As already said, we impose that the matrices M_S and M_x are of full rank m to fully utilize processor resources. There are several cases according to the shape of the matrix F :

$q_x = d$ (in such a case F is square).

This is clearly the simplest case when F is square and non-singular. Given M_x of rank m , then $M_x F$ is of rank m and we can let $M_S = M_x F$ without violating the constraint that M_S is of rank m . Conversely, given M_S of rank m , $M_S F^{-1}$ is of rank m , and we can let $M_x = M_S F^{-1}$.

$q_x < d$ (in such a case, F is flat).

Assume that F is of full rank. Then, given M_x of rank m , $M_x F$ is of rank ² m , and we can safely let $M_S = M_x F$. However, given M_S , finding M_x of rank m such that $M_S = M_x F$ is not always possible. We know that F admits a right pseudo-inverse³ F^{-1} of size $d \times q_x$ and of rank q_x , such that $FF^{-1} = Id$. Hence, if there exists M_x such that $M_S = M_x F$, then $M_S F^{-1} = M_x FF^{-1} = M_x$. Unfortunately, $M_x = M_S F^{-1}$ is not always a solution of the equation $M_S = M_x F$. we have the compatibility condition⁴ $M_S = M_S F^{-1} F$. Furthermore, $M_S F^{-1}$ can be of arbitrary rank less than m . To summarize, given M_x of rank m , it is always possible to determine M_S of rank m ($M_S = M_x F$) while the converse is not true.

$d < q_x$ (in such a case, F is narrow).

Assume that F is of full rank d . Then, F has a left pseudo-inverse F^{-1} of size $q_x \times d$, of rank d , and such that $F^{-1}F = Id$. The situation is exactly the converse of the previous one: given M_S of rank m , $M_x = M_S F^{-1}$ is a rank- m solution⁴ of the equation $M_S - M_x F = 0$; however, given M_x of rank m , $M_x F$ can be of arbitrary rank less than m . Hence, given M_S of rank m , it is always possible to determine M_x of rank m ($M_x = M_S F^{-1}$) while the converse is not true.

2.2.2 Access graph

We recall here the definition of a m -dimensional access graph $G = (V, E, m)$ graph as stated in [6]. The access graph takes into account how the equation $M_S = M_x F$ can be solved.

1. m is the dimension of the target virtual architecture,
2. each vertex $v \in V$ represents an array variable or a statement,
3. consider a loop nest where an array variable x of dimension q_x is accessed (read or written) in a statement S of depth d , through an access matrix F of rank $\min(q_x, d)$ greater than the dimension m of the target architecture: then if $q_x \leq d$ we have an edge from x to S , to indicate that given M_x of rank m it is always possible to find M_S of rank m such that the

²See Lemma 1 in the appendix.

³See the appendix for background on pseudo-inverses.

⁴See Lemma 3 in the appendix.

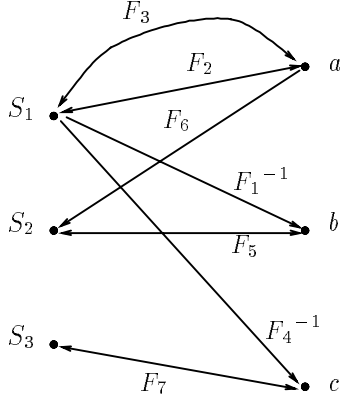


Figure 1: Access graph ($m = 1$ or $m = 2$)

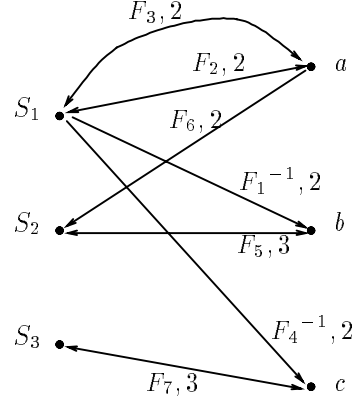


Figure 2: Access graph with integer weights

communication is made local, the weight of the edge is F ; and if $d \leq q_x$ we have an edge from S to x , to indicate that given M_S of rank m it is always possible to find M_x of rank m such that the communication is made local, the weight of the edge is F^{-1} . In the special case where $q_x = d$, there is a single edge with two arrows between S and x to indicate that both orientations are possible.

Figure 1 represents the access graph for our example. The communication corresponding to the access matrix F_8 is not represented because F_8 is not a full rank matrix.

Remark For a narrow $q_x \times d$ matrix F , the left pseudo-inverse is not the only matrix G that satisfies the equation $GF = \text{Id}$. Any matrix H such that

$$H = F^{-1} + M(\text{Id} - FF^{-1})$$

where M is an arbitrary matrix of correct dimension also satisfies $GF = \text{Id}$. Hence, in the access graph, we can choose any full rank matrix G such that $GF = \text{Id}$ as weight matrix (and not necessarily the “true” left pseudo-inverse F^{-1} as defined in [14]).

2.2.3 Mapping heuristic

Some nonlocal communications represented in the graph can be zeroed out, but not all. Consider a simple path in the access graph going from vertex v_1 to vertex v_2 with $v_1 \neq v_2$ (the path is not a cycle). Then given any allocation matrix M_{v_1} of rank m for vertex v_1 , the existence of the path ensures that it is always possible to make local all communications represented by edges between the two vertices. In our example, given M_a of rank 2 in $G = (V, E, 2)$, and following the path $a \rightarrow S_1 \rightarrow b \rightarrow S_2$, we are ensured to be able to compute M_{S_1} , M_b and M_{S_2} , all of rank 2, so that the communications corresponding to access matrices F_2 , F_1 and F_5 are made local: we successively let $M_{S_1} = M_a F_2$, $M_b = M_{S_1} F_1^{-1}$ and $M_{S_2} = M_b F_5$.

There is another path from a to S_2 with the direct edge $a \rightarrow S_2$ (access matrix F_5 for reading a in S_2). Is it possible to make local this communication in addition to the three above communications? Using the edge $a \rightarrow S_2$ we get the equation $M_{S_2} = M_a F_5$ while we had $M_{S_2} = M_b F_5 = M_{S_1} F_1^{-1} F_5 = M_a F_2 F_1^{-1} F_5$.

To simplify the equations, in the following we use the matrix $F'_1 = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 0 & -1 \end{pmatrix}$ (which satisfies the equation $GF_1 = \text{Id}$) instead of F_1^{-1} . We derive the condition $M_a F_2 F'_1 F_5 = M_a F_6$. This condition is satisfied for all matrices M_a of rank m iff $F_2 F'_1 F_5 = F_6$. In our example, let $F_{path} = F_2 F'_1 F_5$. We compute $F_{path} = \begin{pmatrix} 0 & 0 & -1 \\ 1 & 0 & 1 \end{pmatrix} \neq F_6$.

Note that as $F_{path} - F_6 = \begin{pmatrix} 0 & -1 & 0 \\ 0 & 1 & 0 \end{pmatrix}$ is of deficient rank, according to m , it will or not be possible to find a matrix M_a such that the condition $M_a F_2 F'_1 F_5 = M_a F_6$ is satisfied.

In fact, this analysis can be extended in the general case: each time there are two disjoint paths p_1 and p_2 both going from a vertex v_1 to a vertex v_2 in the access graph, we can make all communications on both paths local provided that the equality $F_{p_1} = F_{p_2}$ holds (where F_p denotes the product of the access matrices along the edges of path p). If $F_{p_1} - F_{p_2}$ is of deficient rank, according to the size of the allocation matrix, it can or not be possible to find a matrix M such that $M(F_{p_1} - F_{p_2}) = 0$. See [6] for more details.

Besides, a similar analysis can be performed in the case of cycles. Denote by F_c the product of the weight matrices along the cycle: if $F_c = \text{Id}$ (where Id is the identity matrix), all the communications can be made local along the cycle, if $F_c - \text{Id}$ is of deficient rank, according to the size of the allocation matrix, it can or not be possible to have only local communications along the cycle.

As already said, the access graph depends upon the dimension m of the target architecture space. Given m , not all the communications are taken into account in the access graph $G = (V, E, m)$. The edges in G represent only the communications with access matrix of full rank greater than m . So, the heuristic does not try to make all the communications local but only the “most important ones”. We use the following heuristic:

Heuristic Given the access graph $G(V, E, m)$:

1. associate an integer weight to each edge (see below). Construct a maximum branching $G' = (V', E', m)$ of G using the algorithm due to Edmonds,
2. for each edge in $E \setminus E'$, try to add the edge to G' . If the addition of the new edge creates a cycle of matrix weight the matrix identity or a new path with same source and destination vertices and same weight as an already existing path, the edge can be added in E' . At this step, all the communications represented by edges in G' can be made local,
3. consider the multiple paths and the cycles with $F_{p_1} - F_{p_2}$ or $F_{cycle} - I$ of deficient rank and try to find allocation matrices that allow to zero out even these communications.

In step 1, the weight of a branching is defined as the sum of the integer weights of the arcs in the branching. A maximum branching is any possible branching with the largest possible weight, see [7]. The simplest weight function is to assign the same value 1 to all edges of G . But we can give priority to edges that involve a large volume of communication. We do not need to precisely know the volume of data exchanged, a consistent estimate is sufficient. Consider the following loop nest:

```

for I do
    ...
    { Statement S } a(FaI + ca) = ... b(FbI + cb) ...
endfor

```

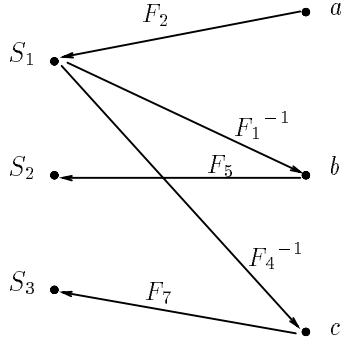


Figure 3: A possible branching

We propose to use the dimension of the set $\{J \mid \exists I, J = F_a I + c_a\}$ (resp $\{J \mid \exists I, J = F_b I + c_b\}$) to estimate the volume of data exchanged for the communication corresponding to the access matrix F_a (resp, F_b). Hence, in the access graph, an integer weight corresponding to the rank of the access matrix is associated to each edge. In this way, communications inducing the largest traffic are zeroed out in priority. Note that in the heuristic proposed by Feautrier [9], the volume of data exchanged induced by a communication is taken into account in a similar way.

Figure 2 represents the weighted access graph for our motivating example. A matrix weight (corresponding to the access matrix or its pseudo-inverse) and an integer weight (corresponding to the volume of data exchange due to the communication) are associated to each edge.

A possible maximum branching for our example is represented in Figure 3. It contains 5 edges out of the 7 edges of the access graph. Hence, 5 communications can be made local and 2 communications remain nonlocal communications. Note that both edges of maximum weight 3 have been zeroed out.

2.3 Optimizing residual communications

The main problem now is: what to do with the residual general communications? After the previous heuristic, the communication corresponding to the access matrix F_6 (reading a in S_1) and the one corresponding to the access matrix F_3 (reading a in S_1 too) remain nonlocal general communications.

To simplify the equations, we use the matrix $F'_1 = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 0 & -1 \end{pmatrix}$ (which satisfies the equation $GF_1 = \text{Id}$) instead of F_1^{-1} and $F'_4 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}$ (which satisfies $GF_4 = \text{Id}$) instead of F_4^{-1} (see the remark in Section 2.2.2).

Back to the branching, there is one input vertex a . Hence, we can choose the allocation matrix M_a and deduce the other allocation matrices from M_a in order to make the communications local: $M_{S_1} = M_a F_2$, $M_b = M_{S_1} F'_1$, $M_c = M_{S_1} F'_4$, $M_{S_2} = M_b F_5$ and $M_{S_3} = M_c F_7$.

Choose, for example, $M_a = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$. We have $M_{S_1} = M_a F_2 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$, $M_b = M_a F_2 F'_1 = \begin{pmatrix} 0 & 0 & -1 \\ 1 & 0 & 1 \end{pmatrix}$, $M_{S_2} = M_a F_2 F'_1 F_5 = \begin{pmatrix} 0 & 0 & -1 \\ 1 & 0 & 1 \end{pmatrix}$, $M_c = M_a F_2 F'_4 = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}$, and $M_{S_3} =$

$$M_a F_2 F_4 F_7 = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}.$$

We check that the two communications corresponding to the two edges of the access graph that do not belong to the selected branching remain nonlocal communications:

$$M_a F_6 = \begin{pmatrix} 0 & 1 & -1 \\ 1 & -1 & 1 \end{pmatrix} \neq M_{S_2} \text{ and } M_a F_3 = \begin{pmatrix} 1 & 1 \\ -4 & -5 \end{pmatrix} \neq M_{S_1}.$$

Remark If we left-multiply M_a by a unimodular matrix U of $M_n(\mathbb{Z})$ ($M_n(\mathbb{Z})$ denotes the set of $n \times n$ matrices over \mathbb{Z} and the unimodular matrices of $M_n(\mathbb{Z})$ are those of determinant ± 1), all the allocation matrices deduced from M_a will be left-multiplied by the same unimodular matrix. Inside each connected component of the branching, the alignment matrices are computed up to a multiplication by an unimodular matrix.

How to optimize the two residual communications?

1. Detecting macro-communications

For the first communication concerning the access matrix F_6 ($a(F_6(i, j)^t + c_6)$ read in S_2), we can notice that F_6 has a non null kernel: $\ker F_6$ is the subset generated by the vector $v = (0, 1, 1)^t$.

Let I_1 and I_2 two points of the index set of S_2 such that $I_2 - I_1 \in \ker F_6$, i.e $I_2 - I_1 = kv$, $k \in \mathbb{Z}$: we have $F_6 I_1 = F_6 I_2$. Besides, $M_{S_2} v = (-1, 1)^t$, hence $M_{S_2} I_1 \neq M_{S_2} I_2$.

The same value of array a , $a(F_6 I_1 + c_6)$ (or $a(F_6 I_2 + c_6)$), located in the memory of processor $M_a(a(F_6 I_1 + c_6)) + \alpha_a$, must be read by distinct processors $M_{S_2} I_1 + \alpha_{S_2}$ and $M_{S_2} I_2 + \alpha_{S_2} = M_{S_2}(I_1 + kv) + \alpha_{S_2}$. Hence, the communication can be viewed as a partial broadcast along one direction of the processor space.

Macro-communications such as broadcasts are efficiently implemented on modern DMPCs. On a CM-5, the ratio between the communication time for a general communication and a broadcast is 60 (see Table 1 in section 1). The broadcast can be total (the value is sent to the whole processor space) or partial (the value is sent in only some directions of the processor space). To be most efficient a partial broadcast must be performed along the directions of the processor space. To optimize the first residual communications, we choose the allocation matrices so as to have a partial broadcast along one axis of the processor space.

In our example, the value $a(F_6 I_1 + c_6)$ in the memory of the processor $M_a(a(F_6 I_1 + c_6)) + \alpha_a$ is sent to processors $M_{S_2}(I_1 + kv) + \alpha_{S_2}$, $k \in \mathbb{Z}$. The communication can be decomposed in a translation and a partial broadcast along the direction given by the vector $M_{S_2} v$. However,

$$M_{S_2} v = \begin{pmatrix} 0 & 0 & -1 \\ 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} -1 \\ 1 \end{pmatrix}.$$

With such a mapping, the broadcast is not parallel to an axis. We *rotate* the mapping by left multiplying M_a (and therefore all the other allocation matrices) by a suitable unimodular matrix. Let $V = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$. We have $V M_{S_2} v = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$. By left-multiplying all the allocations

matrices by the unimodular matrix V , the communication corresponding to the access matrix F_6 becomes a macro-communication that can be performed efficiently⁵.

2. Decomposing the last residual communication

For the second communication concerning the access matrix F_3 ($a(F_3(i, j)^t + c_3)$ read in S_1), we try to decompose it into more simple and efficient data movements such as horizontal or vertical ones.

In our example, the processor $P = VM_a F_3 I$ sends its data to the processor $Q = VM_{S_1} I$ (up to a translation). Let T be the data flow matrix: a processor P sends data to processor $Q = TP$. We have $TVM_a F_3 = VM_{S_1}$. So, $T = VM_{S_1}(M_a F_3)^{-1}V^{-1}$. Besides

$$VM_{S_1}(M_a F_3)^{-1}V^{-1} = \begin{pmatrix} 1 & -1 \\ 5 & -4 \end{pmatrix}.$$

We can decompose $VM_{S_1}(M_a F_3)^{-1}V^{-1}$ into two elementary matrices that correspond to horizontal or vertical communications (see Section 4):

$$VM_{S_1}(M_a F_3)^{-1}V^{-1} = \begin{pmatrix} 1 & 0 \\ 5 & 1 \end{pmatrix} \begin{pmatrix} 1 & -1 \\ 0 & 1 \end{pmatrix}.$$

To summarize, in our example, we finally obtain on the access graph 5 local communications, one broadcast and one residual communication that can be decomposed into two elementary communications.

3 Macro-communications

In this Section we derive formal conditions for detecting and implementing macro-communications such as broadcasts, scatters, gathers, and reductions. We also address the message vectorization problem.

3.1 Broadcast

Broadcasts occur when the same data item is accessed at same time-step by several virtual processors. Consider the following loop nest:

Example 2

for I do

...

S(I) ... = $a(F_a I + c_a)$

...

endfor

⁵We point out that the rank-deficient communication corresponding to F_8 also becomes a broadcast parallel to one direction of the processor space: $\ker F_8$ is the subset generated by the vectors $v_1 = (0, 1, -1)^t$ and $v_2 = (1, 0, -1)^t$. We have $VM_{S_3}[v_1 \ v_2] = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$. Of course this is a lucky coincidence !

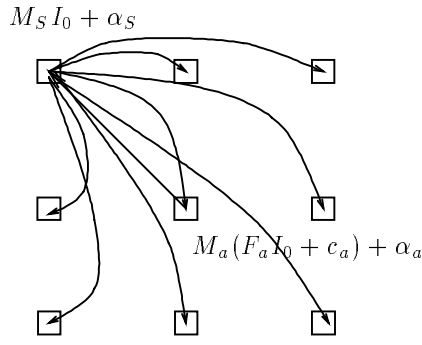


Figure 4: Complete broadcast, $m = 2, p = 2$

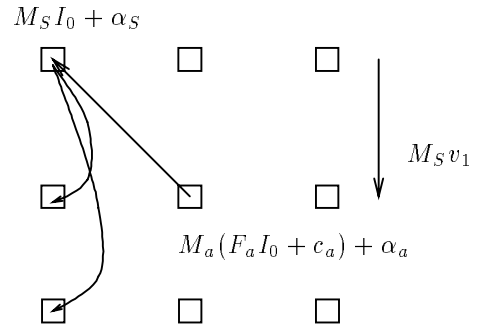


Figure 5: Partial broadcast, $m = 2, p = 1$

Let $M_S I + \alpha_S$ and $M_a I + \alpha_a$ be the affine allocation functions for statement S and for array a . We assume that the computation time steps for $S(I)$ are given by a linear multi-dimensional schedule. Let Θ_S be the multi-dimensional scheduling application for statement S , the computation of $S(I)$ is scheduled at time-step $t = \Theta_S(I)$ (see [8]) on the processor $M_S I + \alpha_S$ and the data accessed is $a(F_a I + c_a)$ which is located in the processor $M_a(F_a I + c_a) + \alpha_a$.

The same index x from array a is read at same time-step by several processors if there exist two indices I_1, I_2 of the iteration space such that:

1. $t = \Theta_S I_1 = \Theta_S I_2$,
2. $x = F_a I_1 + c_a = F_a I_2 + c_a$,
3. $M_S I_1 + \alpha_S \neq M_S I_2 + \alpha_S$.

This implies that $I_1 - I_2 \in \ker(\Theta_S) \cap \ker(F_a) \setminus \ker(M_S)$. Let IS be the index set for statement S . Let $I_0 \in \text{IS}$, let p the dimension of $\ker(\Theta_S) \cap \ker(F_a) \setminus \ker(M_S)$ and (v_1, v_2, \dots, v_p) vectors that generate $\ker(\Theta_S) \cap \ker(F_a) \setminus \ker(M_S)$. Let $v \in \ker(\Theta_S) \cap \ker(F_a) \setminus \ker(M_S)$, the computations of $S(I_0)$ and $S(I_0 + v)$ are scheduled at the same time-step and the same element $a(F_a I_0 + c_a)$ of the array a is read. The (value of the) data item $a(F_a I_0 + c_a)$ is sent to both processors $M_S I_0 + \alpha_S$ and $M_S(I_0 + v) + \alpha_S$. Hence, the necessary communications for the set of computations $\{S(I) \mid I = I_0 + v, v \in \ker(\Theta_S) \cap \ker(F_a) \setminus \ker(M_S)\}$ can be regrouped into two communications: first a translation of the data item $a(F_a I_0 + c_a)$ from $M_a(F_a I_0 + c_a) + \alpha_a$ to $M_S I_0 + \alpha_S$, then a broadcast of this item along the vectors $M_S v_1, M_S v_2, \dots, M_S v_p$.

Let m be the dimension of the target virtual processor space. If $p = m$, the broadcast is total. If $0 < p < m$, the broadcast is partial. See Figures 4 and 5 for examples with $m = 2$ and $p = 2$ or $p = 1$. If $p = 0$, the broadcast is hidden by the mapping and we have only a point to point communication. In the case of a partial broadcast, as Platonoff in [19], we impose to broadcast in directions parallel to some axis of the processor space. In Figure 5, the direction of the broadcast $M_S v_1$ is parallel to one axis.

Partial broadcast conditions Let D be the $m \times p$ matrix $\begin{bmatrix} M_S v_1 & M_S v_2 & \dots & M_S v_p \end{bmatrix}$. In the case of partial broadcast, $0 < p < m$, D is a narrow rectangular matrix. Partial broadcasts correspond to efficient schemes of communications iff there are along some dimensions of the processor space, i.e., $D = \begin{bmatrix} D_1 \\ 0 \end{bmatrix}$ (up to a row permutation), where D_1 is $n \times p$ full rank matrix ($n \leq p$) and 0 is a $(m - n) \times p$ null matrix.

If D is not of the previous form, we use the right Hermite form of D ⁶ to find a new allocation matrix M_S such that the broadcast is made parallel to some axis: we decompose D as $D = Q \begin{bmatrix} H \\ 0 \end{bmatrix}$ where Q is a unimodular matrix and H is $n \times n$ lower triangular matrix. Hence, $\begin{bmatrix} Q^{-1}M_S v_1 & Q^{-1}M_S v_2 & \dots & Q^{-1}M_S v_p \end{bmatrix} = \begin{bmatrix} H \\ 0 \end{bmatrix}$. If we left multiply M_S by the unimodular matrix Q^{-1} , the partial broadcast becomes parallel to the axis.

3.2 Scatter

A scatter occurs when several data located in the same processor must be sent at same time-step to several processors. The only difference between broadcast and scatter is that different data is to be sent to the receiving processors. Consider again Example 2, the conditions to have a scatter are the following:

1. $t = \Theta_S I_1 = \Theta_S I_2$,
2. $p = M_a(F_a I_1 + c_a) + \alpha_a = M_a(F_a I_2 + c_a) + \alpha_a$,
3. $M_S I_1 + \alpha_S \neq M_S I_2 + \alpha_S$,
4. $F_a I_1 + c_a \neq F_a I_2 + c_a$.

This implies that $I_1 - I_2 \in (\ker \Theta_S \cap \ker M_a F_a) \setminus (\ker M_S \cap \ker F_a)$. We have similar conditions for a scatter as for a broadcast.

3.3 Gather

A gather is the “inverse” operation to a scatter. Several data located in different processor are sent at same time-step to the same processor. Consider the following loop nest:

Example 3

for I do

 ...

 S(I) $a(F_a I + c_a) = \dots$

 ...

endfor

Again, the conditions for a gather are close to that for scatters:

1. $t = \Theta_S I_1 = \Theta_S I_2$,
2. $p = M_a(F_a I_1 + c_a) + \alpha_a = M_a(F_a I_2 + c_a) + \alpha_a$,
3. $M_S I_1 + \alpha_S \neq M_S I_2 + \alpha_S$,
4. $F_a I_1 + c_a \neq F_a I_2 + c_a$.

This implies that $I_1 - I_2 \in (\ker \Theta_S \cap \ker M_a F_a) \setminus (\ker M_S \cap \ker F_a)$. A gather can be partial or total and we have similar conditions on partial gathers as previously.

⁶See Definition 1 in the appendix

3.4 Reduction

A reduction is similar to a gather but the different values sent to the same processor are used to compute one single value. Reductions occur when, at same time-step, a single processor uses values computed by different instances of the same instruction on different processors. A reduction is usually associated with a commutative and associative function $(+, \min, \dots)$ that computes a resulting value from many input values. Consider the following loop nest, where s represents an array element:

Example 4

```
for  $I$  do
  ...
  S(I)  $s = s + b(F_b I + c_b)$ 
  ...
endfor
```

The conditions to have a reduction are:

1. $t = \Theta_S I_1 = \Theta_S I_2$,
2. $M_b(F_b I_1 + c_b) + \alpha_b \neq M_b(F_b I_2 + c_b) + \alpha_b$,
3. $M_S I_1 = M_S I_2$.

This implies that $I_1 - I_2 \in \ker \Theta_S \cap \ker M_S \setminus \ker M_b F_b$.

3.5 Message vectorization

Message vectorization can take place when a processor accesses data from another processor that remains unchanged for several consecutive time steps: data items to be sent can be regrouped into packets that are sent just in time to reach their destination. The communications can be extracted out of a loop and performed before the computations. The idea is to replace a set of small-size communications by a single large message, so as to reduce overhead due to start-up and latency [10, 21, 5].

Consider again Example 2. The space-time transformation is given by:

$$\begin{pmatrix} t \\ p \end{pmatrix} = \begin{pmatrix} \Theta_S \\ M_S \end{pmatrix} I + \begin{pmatrix} 0 \\ \alpha_S \end{pmatrix},$$

where p is the processor responsible at time-step t (the schedule can be multi-dimensional) for the computation $S(I)$. Let $\tilde{S} = \begin{pmatrix} \Theta_S \\ M_S \end{pmatrix}$. Hence, we have $I = \tilde{S}^{-1} \begin{pmatrix} t \\ p - \alpha_S \end{pmatrix}$. For the computation $S(I)$, processor p reads data from processor $M_a(F_a I + c_a) + \alpha_a$. This expression does not depend on t when $M_a F_a \tilde{S}^{-1} = (0, X)$ where 0 is a $m \times p$ null matrix and X is a $m \times m$ matrix. This condition is equivalent to $M_a F_a = (0, X) \tilde{S} = X M_S$ and therefore to $\ker M_S \subset \ker(M_a F_a)$.

Communication	Not decomposed	L	U	$L.U$
Execution ratio	10	1.5	3.9	5.4

Table 2: Decomposing versus not decomposing a general affine communication on the Intel Paragon

of the larger grid dimension. Data was distributed using a standard `CYCLIC` distribution. We come back to distribution techniques in Section 4.3 where we introduce a new data distribution scheme that enables to implement elementary communications even more efficiently.

Anyway, the gain obtained by decomposing communication is machine- and compiler-dependent. Table 2 is only intended to give experimental evidence that communication decomposition can prove efficient. To be conservative, we look to decomposing general communications into a small (say $l \leq 4$) number of elementary ones.

4.2 Communication decomposition: how

4.2.1 Direct decomposition

Let $T = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$ with $\det(T) = ad - bc = 1$. We derive conditions to know whether T can be decomposed into the product or less than or equal to four elementary matrices L_i and U_i .

The product of two elementary matrices is

$$LU = \begin{pmatrix} 1 & 0 \\ \beta & 1 \end{pmatrix} \times \begin{pmatrix} 1 & \alpha \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & \alpha \\ \beta & 1 + \alpha.\beta \end{pmatrix}, \text{ or } UL = \begin{pmatrix} 1 & \alpha \\ 0 & 1 \end{pmatrix} \times \begin{pmatrix} 1 & 0 \\ \beta & 1 \end{pmatrix} = \begin{pmatrix} 1 + \alpha.\beta & \alpha \\ \beta & 1 \end{pmatrix}$$

hence the necessary and sufficient condition to have such a decomposition is: $a = 1$ or $d = 1$.

We see that the conditions on (a, b, c, d) to have a decomposition beginning with a U are the same as the conditions on (d, c, b, a) for a decomposition beginning with a L . Thus we can restrict ourselves to decompositions beginning with a U .

The equation for a decomposition into the product of three elementary matrices (beginning with a U matrix) is:

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{pmatrix} 1 & \alpha \\ 0 & 1 \end{pmatrix} \times \begin{pmatrix} 1 & 0 \\ \beta & 1 \end{pmatrix} \times \begin{pmatrix} 1 & 0 \\ \gamma & 1 \end{pmatrix} = \begin{pmatrix} 1 + \alpha.\beta & \alpha + \beta(1 + \beta.\gamma) \\ \beta & 1 + \beta.\gamma \end{pmatrix}$$

So $a = 1 + \alpha.\beta$, $b = \alpha + \beta(1 + \beta.\gamma)$, $c = \beta$, and $d = 1 + \beta.\gamma$. We have $\alpha = b - \gamma.d = b - c.d$, and $a = 1 + \alpha.\beta = (b - c.d).c$. Thus c needs to divide $a - 1$. As $\det(T) = 1$, this condition is sufficient too. To summarize, the necessary and sufficient condition for T to be equal to the product of three elementary matrices is: c divides $a - 1$ or b divides $d - 1$.

The equation for a decomposition into the product of four elementary matrices (beginning with a U matrix) is:

$$\begin{pmatrix} 1 & \alpha \\ 0 & 1 \end{pmatrix} \times \begin{pmatrix} 1 & 0 \\ \beta & 1 \end{pmatrix} \times \begin{pmatrix} 1 & \gamma \\ 0 & 1 \end{pmatrix} \times \begin{pmatrix} 1 & 0 \\ \delta & 1 \end{pmatrix} = \begin{pmatrix} (1 + \alpha.\beta).(1 + \gamma.\delta) + \alpha.\delta & (1 + \alpha.\beta).\gamma + \alpha \\ \beta.(1 + \gamma.\delta) + \delta & 1 + \beta.\gamma \end{pmatrix}$$

Thus $1 + \beta.\gamma = d$. Then

$$(1 + \alpha.\beta).\gamma = b \Rightarrow \alpha.(1 + \beta.\gamma) = b, \text{ so } \alpha.\delta + \gamma = b, \gamma = b - \alpha.d \quad (1)$$

$$\beta.(1 + \gamma.\delta) + \delta = c \Rightarrow \delta.(1 + \beta.\gamma) + \beta = c, \text{ so } \delta.d + \beta = c, \beta = c - \delta.d \quad (2)$$

$(1 + \alpha.\beta)(1 + \gamma.\delta) + \alpha.\delta = a$ is always true. Indeed

$$(1 + \alpha.\beta)(1 + \gamma.\delta) + \alpha.\delta = 1 + \alpha.\beta + \gamma.\delta + \alpha.\beta.\gamma.\delta + \alpha.\delta \quad (3)$$

$$= 1 + \alpha.(c - \delta.d) + (b - \alpha.d).\delta + \alpha.\delta(d - 1) + \alpha.\delta \quad (4)$$

$$= 1 + \alpha.c + \delta.b - \alpha.\delta.d = a \quad (5)$$

As $d = 1 + \beta.\gamma = 1 + (c - \delta.d).(b - \alpha.d)$. Thus $d - 1 = b.c - \delta.d.b - \alpha.d.c + \alpha.\delta.d^2$, with $bc = ad - 1$ because $\det(T) = 1$. So the equation can be divided by d : $1 = a - \alpha.c - \delta.b + \alpha.\delta.d$ so $a - 1 - \alpha.c = \delta.(b - \alpha.d)$. Accordingly an integer α needs to exist such that $b - \alpha.d$ divides $a - 1 - \alpha.c$.

Given $r = b - \alpha.d$, then r divides $a - 1 - \alpha.c$, or

$$a - 1 - \alpha.c = a - 1 - c.(b - r)/d \quad (6)$$

$$= (a.d - d - b.c - r.c)/d \quad (7)$$

$$= (1 - d + r.c)/d \quad (8)$$

But $r \wedge d = 1^7$ (indeed if $p|r$ and $p|d$ then $p|b$, and so $p = 1$ because $b \wedge d = 1$). Thus, $r|(1 - d + rc)$, so $r|(1 - d)$. Therefore, the necessary and sufficient condition for a matrix to be equal to the product of four elementary matrices is: $\exists \alpha, b - \alpha.d | d - 1$ or $\exists \alpha, c - \alpha.a | a - 1$.

We could go further and look for a decomposition into five or more elementary matrices. But an exhaustive search shows that every 2×2 matrix T with $\det(T) = 1$ and whose coefficients are all lower than or equal to 14 in absolute value, is equal to the product of 2, 3 or 4 elementary matrices. In practice, larger coefficients are unlikely to be encountered in loop nests !

4.2.2 With left-multiplication by a unimodular matrix

As outlined in Section 2.3, alignment matrices are computed up to a multiplication by an unimodular matrix. If we left-multiply M_S and M_a by a unimodular matrix M , then the dataflow matrix $T = M_S(M_a.F_a)^{-1}$ is transformed into MTM^{-1} . Therefore, rather than decomposing T into the product of elementary matrices, we can search for a unimodular matrix M such that MTM^{-1} , a matrix *similar* to T , can be decomposed into such a product.

Consider again the case $m = 2$ and $\det(T) = 1$. The best would be to show that T is similar to a product LU of two elementary matrices, so as to decompose MTM^{-1} into one horizontal communication followed by one vertical communication. The problem amounts to show that every integer matrix $T = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$ with $ad - bc = 1$ is similar over \mathbb{Z} to a matrix $LU = \begin{pmatrix} 1 & \alpha \\ \beta & 1 + \alpha\beta \end{pmatrix}$. The problem is surprisingly difficult, and the answer depends upon the coefficients of T .

Let $P(X) = X^2 - (a + d)X + 1$ be irreducible over \mathbb{Q} : it is the case if $\sqrt{Tr(T)^2 - 4} \notin \mathbb{Q}$, which reduces to $Tr(T) \neq \pm 2$. Then there is a one-to-one correspondence between similarity classes over \mathbb{Z} with the ideal classes in the ring $\mathbb{Z}[\theta]$, where θ is a root of $P(X)$: this is Latimer and MacDuffee's theorem [18, page 53]. In turn, the number of ideal classes is in one-to-one correspondence with the

⁷ $\gcd(a, b)$ is written $a \wedge b$

number of equivalence classes of the quadratic forms with discriminant $\Delta = Tr(T)^2 - 4$ [11, page 192]. This number is finite (it is called the genus) and can be calculated with complex analytic methods, as explained in [2, page 371].

The answer to our problem is negative in the following case: if $Tr(T) - 2 = a + d - 2 = \alpha\beta$ is a prime number, there are only 4 possibilities for α : $\alpha = \pm 1$ or $\alpha = \pm Tr(T) - 2$. Therefore the product of two elementary matrices belong to at most 4 different similarity classes. There are infinitely many values of $\Delta = Tr(T)^2 - 4$ whose genus is greater than 4, hence infinitely many cases where the answer is negative: an arbitrary integer matrix T with $\det(T) = 1$ is not always similar to the product of two elementary matrices.

To be more concrete, let us give an example where the answer is negative. Consider the following example where $Tr(T) = -2$: $T = \begin{pmatrix} 6 & 7 \\ -7 & -8 \end{pmatrix}$. This case is simple because $P(X) = X^2 + 2X + 1$ is not irreducible over \mathbb{Q} . Since $Tr(A) = -2$, we have $\beta\alpha = -4$. Assume there exists a unimodular matrix $M = \begin{pmatrix} u & v \\ w & x \end{pmatrix}$, with $\det(M) = ux - vw = \pm 1$, such that $MTM^{-1} = \begin{pmatrix} 1 & \alpha \\ \beta & 1 + \alpha\beta \end{pmatrix}$. Equating coefficients, we obtain $\beta = \pm(7x^2 - 14xw + 7w^2)$. So 7 needs to divide β , and it is impossible.

However, a simple sufficient condition for T to be similar to a product LU or UL can be determined as follows. Let $T = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$ to be similar to $\begin{pmatrix} 1 & \alpha \\ \beta & 1 + \alpha\beta \end{pmatrix}$. Let $e_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ and $f(e_1) = \begin{pmatrix} a \\ c \end{pmatrix}$. So $f(e_1) - e_1 = \begin{pmatrix} a-1 \\ c \end{pmatrix}$. If $c|a-1$, i.e. $a-1 = \lambda c$, choose $e_2 = \begin{pmatrix} \lambda \\ 1 \end{pmatrix}$. Then (e_1, e_2) is a unimodular basis of determinant 1 and $f(e_1) - e_1 = \lambda e_2$. In the new basis (e_1, e_2) , $T = \begin{pmatrix} 1 & a \\ c & \lambda \end{pmatrix}$. Hence T is similar to the product of two elementary matrices. Our sufficient condition is the same as the necessary and sufficient condition for T to be decomposed into the product of three elementary matrices. Either strategy could be more interesting, depending upon the target machine. Note that all integer matrices T with $\det(T) = 1$ and whose coefficients are all lower than or equal to 5 are similar to a product of 2 elementary matrices.

4.3 A new data distribution scheme

We have dealt with virtual processors so far. In this Section we make a short digression to introduce a new data distribution scheme called *grouped partition*: the grouped partition is well-suited to implementing horizontal/vertical communications on the Paragon. The grouped partition leads to smaller communication times than the standard CYCLIC or CYCLIC(BLOCK) distributions.

Grouped partition Consider a horizontal communication of dataflow matrix $U = \begin{pmatrix} 1 & k \\ 0 & 1 \end{pmatrix}$. Virtual processor (i, j) sends data to processor $(i + kj, j)$. So processor $(0, j)$ sends to (kj, j) which in turn sends to $(2k, j)$ and so on. The communication is thus subdivided into k independent communications. Let $0 \leq c < k$ design a class defined as follows: virtual processor $(c + \alpha k, j)$, $\forall \alpha$ belong to class c . There is no communication between classes, communication occurs only within classes.

Targeting a 2d-grid of $P \times Q$ processors, we have P physical processors per row and Q per column. As P is not expected to be equal to k , physical processors are partitioned into blocks, and

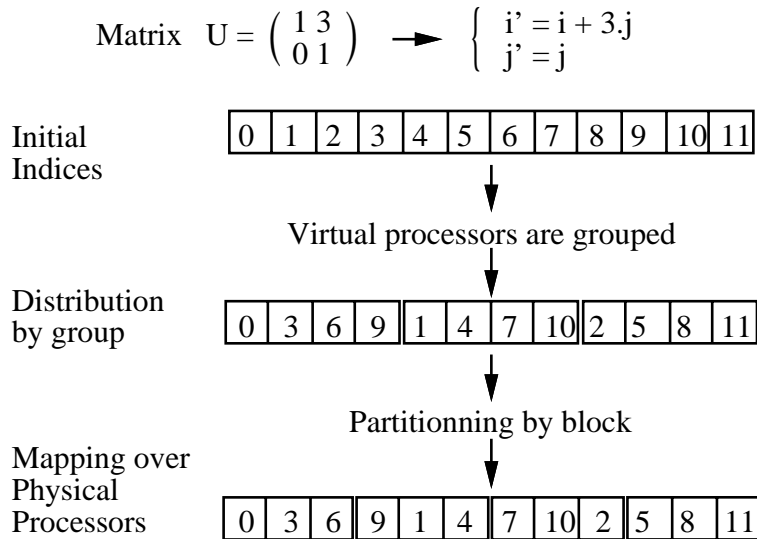


Figure 6: U is the communication matrix. The 12 virtual processors on each row are grouped into 3 classes and mapped onto $P = 4$ physical processors.

blocks are assigned to a class, as illustrated in Figure 6.

Partitioning with a dataflow matrix L is similar in the vertical dimension. This we obtain a two dimensional partition as shown in Figure 7 that is well-suited to implement a product LU or UL (note that communication L and U are performed one after the other, not in parallel).

4.4 Comparison

Experiments have been done on the Intel Paragon in order to compare communication time for the grouped partition with the standard `CYCLIC` and `CYCLIC(BLOCK)` distributions. Figure 8 show the results obtained on a Paragon.

We see that the grouped partition is always more efficient than a standard `BLOCK` or `CYCLIC(BLOCK)` distribution. The `CYCLIC` distribution performs well, because it amounts to the grouped partition with $k = 1$.

We point out that the performance of a general affine communication (not decomposed) is not affected by the new data distribution scheme. Running several experiments with several grid sizes, we observe less than 5% difference between the grouped partition and the `CYCLIC` distribution. Therefore, our new data distribution scheme can be seen as an interesting alternative to standard schemes in the context of loop nest parallelization.

$$\text{Matrix } T = \begin{pmatrix} 1 & 2 \\ 3 & 7 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} 1 & 3 \\ 0 & 1 \end{pmatrix} = L.U \quad \text{which implies: } \begin{cases} i' = i + 3j \\ j' = j \end{cases} \quad \text{then: } \begin{cases} i'' = i' \\ j'' = j' + 2i \end{cases}$$

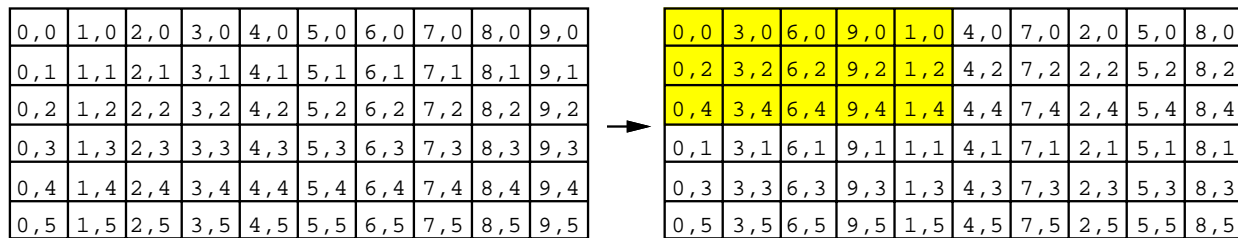


Figure 7: A 10×6 grid of virtual processors is mapped onto a 2×2 grid of physical processors. The dataflow matrix is $T = L.U$. Processor $(0,0)$ is shaded to show the distribution over physical processors.

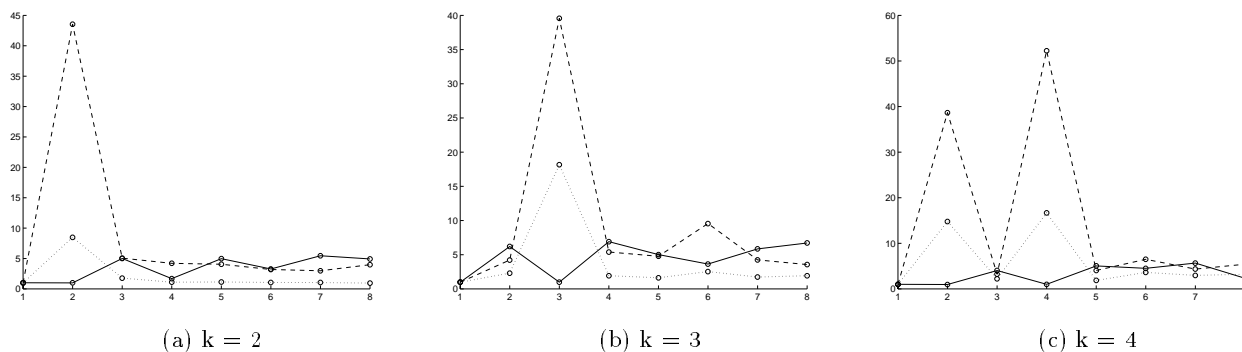


Figure 8: These graphs represent the ratio of the communication time for a U matrix with standard distribution schemes over the time using the grouped partition scheme. For each graph, the dotted line is the ratio of the **CYCLIC(20)** distribution over the grouped partition, the dashed line is the ratio of the full **BLOCK** distribution over the grouped partition, the solid line is the ratio of the **CYCLIC** distribution over the grouped partition.

4.5 Extensions

We have only dealt with dataflow matrices of determinant 1. To generalize to arbitrary matrices, we use a decomposition into “uni-row” or “uni-column” matrices like

$$\begin{pmatrix} 1 & & & & 0 \\ & \ddots & & & \\ \alpha_{0.k} & \cdots & \alpha_{i.k} & \cdots & \alpha_{n.k} \\ & & & \ddots & 0 \\ & & & & 0 & 1 \end{pmatrix}$$

Note that the grouped partition can be used to implement such communications very efficiently.

5 Summary

Summarizing previous Sections, we can sketch our complete heuristic.

1. Zero out non local communications

(a) Access graph

Construct the access graph $G = (V, E, m)$ associated to the loop nest.

(b) Branching

Extract a maximum branching $G' = (V', E', m)$ from the access graph.

(c) Multiple paths, cycles

- i. For each edge in $E \setminus E'$, try to add the edge to G' . If the addition of the new edge creates a cycle of matrix weight the matrix identity or a new path with same source and destination vertices and same weight as an already existing path, the edge can be added in E' . At this step, all the communications represented by edges in G' can be made local.
- ii. In each connected component, consider the multiple paths and the cycles with $F_{p_1} - F_{p_2}$ or $F_{cycle} - I$ of deficient rank and try to find allocation matrices that allow to zero out even these communications.

2. Optimize residual communications

For each connected component of the graph obtained after step 1:

(a) Macro communications

Detect the possible macro-communications and compute the conditions on allocation matrices to have efficient communications after mapping on the virtual processor space. If the conditions are not satisfied, it is possible in each connected component to left-multiply all the allocation matrices by a unimodular matrix (see Section 3).

(b) Decompose residual general communications

Decompose the residual general communications in more simple and efficient ones. If the allocation matrices are not yet fixed in a connected component, to obtain a better decomposition, again left-multiply the allocation matrices by a unimodular matrix (see Section 4).

In the next Section, we present a short review of related work. Then we compare our heuristic to that of Platonoff [19]

6 Related work

6.1 Literature overview

The data alignment problem has motivated a vast amount of research. A brief review of some related work is presented here. The review contains only short abstracts of the presented papers and is by no mean intended to be comprehensive.

Knobe, Lukas and Steele In [13], the authors discuss techniques for automatic layouts of arrays in a compiler targeted to SIMD architectures. The approach to data locality is to consider each occurrence of a datum as a separately allocated object and to mark preferences among certain occurrences to indicate that they should be allocated together. This approach is extended in [17] to MIMD systems. In [16], Lukas shows that same data optimization alignment techniques can be used in both distributed and shared memory systems. For shared memory systems, when alignment preferences can be satisfied, synchronization requirements are eliminated.

Huang and Sadayappan In [12], the authors consider the issue of communication-free hyper-plane partitioning. By modeling the iteration and data spaces and the relation that maps one to another, necessary and sufficient conditions for the feasibility of communication-free partitioning along hyperplanes are characterized.

Li and Chen In [15] the authors formulate the problem of *index domain alignment* as finding suitable alignment functions that embed the index domains of the arrays into a common index domain. The paper contains an NP-completeness result: minimizing the cost of data movement is shown to be an NP-complete problem. Besides, a greedy heuristic algorithm is given at the end of the paper.

Anderson and Lam In [1], the authors propose an algorithm and heuristics that determine alignment for both data and computations (extension of the *owner computes rule*). The algorithm is based on the mathematical model of decompositions as affine functions and is structured into three components: partition, orientation and displacement. The only parallelism exploited is *forall* parallelism or *doacross* parallelism using tiling.

Darte and Robert In [4, 5], the authors introduce a *communication graph* that contains all the information to align data and computations. They formulate ways to reduce the amount of communications (communication rank minimization, broadcasting, message vectorization...). But, the main result is a NP-completeness result. Darte and Robert restrict themselves to a simple case - perfect loop nest in which all access functions are translations - and they show that, even in this case, the alignment problem is NP-complete. They give several heuristics.

Feautrier In [9], Feautrier proposes a greedy algorithm analogous to Gaussian elimination to determine a placement function. Data and computations are aligned in such a way that the *owner computes rule* is respected. The main idea is to zero out edges corresponding to the most important communication volume. An heuristic is given to estimate the communication volume associated to an edge.

Platonoff Platonoff [19] has first raised the following question: as it is usually impossible to zero out (or to make local) all communications, how to efficiently implement the communications that cannot be zeroed out? Platonoff proposed an heuristic based on the greedy heuristic given by Feautrier [9], enlarged by a detection of macro-communications (broadcasts and reductions), whose cost is an order of magnitude smaller than for a general communication (see Table 1).

Platonoff’s algorithm is divided on 4 steps:

1. The rank of the mapping function is defined for each instruction⁸. Then a prototype mapping can be written for each instruction. The prototype mapping is an affine function depending on the the loop indices and the structure parameters⁹.
2. Broadcasts are located in the initial code. To find them, Platonoff uses the *data flow graph*, defined by Feautrier. A broadcast exists if several processors read the same variable at the same time. This means that in a parallel loop, several executions need the same variable. This corresponds to non-zero kernels for data flow graph transformations. A broadcast can be total (one processor sends its data to every other processors) or partial (the data is sent to a subset of the processors).
3. Broadcasts are processed. Broadcast directions are identified in the initial program. The conditions for the prototype mapping are rewritten so that the projection onto the virtual processor space is not along the broadcast directions.
 - (a) Total broadcast conditions: Platonoff first tries to preserve total broadcasts when projecting onto the virtual processor space.
 - (b) Partial broadcast conditions: if a total broadcast is not possible (because it is not compatible with the mapping function rank calculated in the first step), the prototype mapping is augmented with the constraints for a partial broadcast. Another type of constraint is added: to be efficient, partial broadcasts have to be parallel to the axes.
4. The volume of remaining communications is minimized with the greedy heuristic of Feautrier.

6.2 Discussion

Many authors have proposed heuristics to find a communication-free mapping or to the minimize the number of communications. NP-completeness results show that the problem is difficult [15, 4, 6]. A communication-free mapping algorithm usually leads to a trivial mapping (all the computations and the data are regrouped on the same processor) and, therefore, all the parallelism present inside the initial code is lost.

We compare our heuristic with the strategy developed by Platonoff on a small example. Consider the following loop nest:

Example 5

```

for  $t = 1, n$  do
  ...
  for  $i, j, k = 1, n$  do
     $S(I) : a(t, i, j, k) = b(t, i, j)$ 
  
```

⁸The rank is equal the depth of the instruction in the loop nest minus the scheduling dimension.

⁹The structure parameters are integer variables defined outside the loop by an assignment, or by an I/O instruction.

```

    endfor
    ...
endfor

```

Let $\{\vec{e}_i\}_{i=1,4}$ be the canonical basis of the iteration space. We assume here that the loop nest is scheduled by a linear scheduling vector $\pi = e_1^T$. The outer loop is sequential and the inner loops on i, j , and k are parallel. The arrays a and b can be accessed before and after the parallel loop, inside the sequential loop. We also assume that we want to map computations and data onto a

2-dimensional processor space: $m = 2$. The access matrix for array b is $F_b = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}$.

A necessary condition to have a broadcast is $\ker(\pi^t) \cap \ker(F_b) \neq 0$. In our example, we have $\ker(\pi^t) \cap \ker(F_b) = e_4^T$. According to the allocation matrices, the broadcast can be kept or masked. Platonoff's strategy consists in:

1. detecting all the macro-communications possibly present in the initial program,
2. writing the conditions to preserve these macro-communications,
3. making local (or zeroing out) as many remaining communications as possible, by taking into account the conditions previously found in his prototype mapping function.

To keep the partial broadcast, Platonoff would choose $M_S = M_a = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$ and $M_b = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$. With this mapping, each processor sends a value to a column of processors. At each time-step, n^2 broadcasts of one element along one dimension of the processor space are necessary. The loop nest is computable with n^3 macro-communications.

However, we easily see that, with our strategy, the loop nest can be computed without any communication! For example, we can choose $M_b = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$ and $M_S = M_a = M_b F_b = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}$. We first try to make local as many communications as possible and then we try to extract macro-communications from the residual communications, whereas Platonoff first detects the macro-communications and then try to zero out the remaining communications.

7 Conclusion

Many authors have proposed heuristics to minimize the communication volume or number when mapping data and computations of an affine loop nests onto DMPCs. It is generally impossible to obtain a communication-free mapping and another goal in the mapping process is to "optimize" in some sense the residual communications.

We have designed an efficient two-step heuristic

1. based upon the access graph to zero out as many communications as possible, with priority given to communications of largest volume
2. enlarged with the processing of residual communications, either through the extraction of macro-communications or through the decomposition of complex communications into simpler ones

We have provided a detailed analysis of macro-communications (broadcasts, scatters, gathers, reductions) and of message vectorization), together with criteria for their efficient mapping.

Finally we have given analytical formulae to decompose complex communications, and we have shown that such a decomposition improves communication performance on the Paragon.

8 Appendix

8.1 Hermite form

$M_n(\mathbb{Z})$ denotes the set of $n \times n$ matrices over \mathbb{Z} . The unimodular matrices of $M_n(\mathbb{Z})$ are those of determinant ± 1 .

Definition 1 (Right Hermite form) *For every non singular matrix A of $M_n(\mathbb{Z})$, there exist a unimodular matrix $Q \in M_n(\mathbb{Z})$ and a lower triangular matrix $H \in M_n(\mathbb{Z})$ such that:*

- $\forall(i, j), h_{ij} \geq 0$,
- each non-diagonal element is smaller than the diagonal element of the same row,
- $A = QH$.

*This decomposition is called **Hermite form decomposition**. Furthermore, if $\det(A) \neq 0$, Q and H are unique.*

Left Hermite can also be defined in the same way. Besides, Hermite forms also exist for rectangular matrices. Let A be a $(n + t) \times n$ matrix, we can decompose A in the following way:

$$A = Q \begin{bmatrix} H \\ 0 \end{bmatrix}, \text{ where } 0 \text{ is } t \times n \text{ null matrix and } H \text{ is a } n \times n \text{ matrix of the previous form.}$$

8.2 Pseudo-inverses

Let X be a rectangular $u \times v$ integer matrix, and assume that X is of full rank $\min(u, v)$. If $u = v$, then X is nonsingular and its inverse matrix X^{-1} is such that $XX^{-1} = X^{-1}X = Id_u$, where Id_u denotes the identity matrix of order u .

If $u \neq v$, we can define a pseudo-inverse (still denoted as X^{-1}) as follows:

- if $u \leq v$ (X is flat), then XX^t is a square nonsingular $u \times u$ matrix whose (ordinary) inverse matrix is $(XX^t)^{-1}$. Then we define the pseudo-inverse (or right-inverse) of X as $X^{-1} = X^t(XX^t)^{-1}$: X^{-1} is a $v \times u$ matrix of rank u such that $XX^{-1} = Id_u$. Note that $X^{-1}X \neq Id_v$ if $u \neq v$.
- if $u \geq v$ (X is narrow), then X^tX is a square nonsingular $v \times v$ matrix whose (ordinary) inverse matrix is $(X^tX)^{-1}$. Then we define the pseudo-inverse (or left-inverse) of X as $X^{-1} = (X^tX)^{-1}X^t$: X^{-1} is a $v \times u$ matrix of rank v such that $X^{-1}X = Id_v$. Note that in general $XX^{-1} \neq Id_u$ if $u \neq v$.

Note that for square non singular matrices, the pseudo-inverse matrix coincides with the (usual) inverse matrix. For more details, see [14].

8.3 Matrix equations

Lemma 1 *Let A be a $m \times a$ matrix of rank m and F be a $a \times d$ matrix of rank a , where $m \leq a \leq d$. Then AF is of rank m .*

Proof We use the Hermite normal form of F : $F = [H, 0]Q$, where H is a $a \times a$ upper triangular matrix of rank a , and Q is a unimodular $d \times d$ matrix. Since Q is nonsingular, the rank of AF is that of $A[H, 0]$, hence that of AH , hence finally that of A , as H is nonsingular too. ■

Lemma 1 was used in Section 2 to prove that we can safely let $M_S = M_x F$ when M_x is a $m \times q_x$ matrix of rank m and F a $q_x \times d$ matrix of rank q_x , where $m \leq q_x \leq d$. Now for the case where $m \leq d \leq q_x$, we need to solve the equation $M_S = M_x F$, where M_S of rank m and F of rank d are given. We use the following result from [14]:

Lemma 2 *Let S be a $m \times d$ matrix of rank m and F be a $a \times d$ matrix of rank d . Then the equation $XF = S$ admits a solution if and only if the compatibility condition $SF^{-1}F = S$ is satisfied. In such a case, all solutions are given by the expression $X = SF^{-1} + Y(Id_a - FF^{-1})$, where Y is an arbitrary $m \times a$ matrix.*

Lemma 3 *Let S be a $m \times d$ matrix of rank m and F be a $a \times d$ matrix of rank d , where $m \leq d \leq a$. Then the equation $XF = S$ admits the rank- m solution $A = SF^{-1}$.*

Proof The compatibility condition is verified because $F^{-1}F = Id_d$ with our hypothesis. Hence $A = SF^{-1}$ is a solution of the equation. Finally, we apply Lemma 1 to prove that its rank is indeed m . ■

Lemma 3 was used in Section 2 to orient some arrows from statements to arrays in the access graph.

References

- [1] Jennifer M. Anderson and Monica S. Lam. Global optimizations for parallelism and locality on scalable parallel machines. *ACM Sigplan Notices*, 28(6):112–125, June 1993.
- [2] J. W. S. Cassels. *Rational Quadratic Forms*. Academic Press, 1978.
- [3] Tzung-Shi Chen and Jang-Ping Sheu. Communication-free data allocation techniques for parallelizing compilers on multicomputers. *IEEE Transactions on Parallel and Distributed Systems*, 5(9):924–938, 1994.
- [4] Alain Darte and Yves Robert. Mapping uniform loop nests onto distributed memory architectures. *Parallel Computing*, 20:679–710, 1994.
- [5] Alain Darte and Yves Robert. On the alignment problem. *Parallel Processing Letters*, 4(3):259–270, 1994.
- [6] Michèle Dion and Yves Robert. Mapping affine loop nests: New results. In Bob Hertzberger and Guiseppe Serazzi, editors, *High-Performance Computing and Networking, International Conference and Exhibition*, volume LCNS 919, pages 184–189. Springer-Verlag, 1995. Extended version available as Technical Report 94-30, LIP, ENS Lyon (anonymous ftp to lip.ens-lyon.fr).

- [7] J.R. Evans and E. Minieka. *Optimization Algorithms for Networks and Graphs*. Marcel Dekker Inc, 1992.
- [8] Paul Feautrier. Some efficient solutions to the affine scheduling problem, part II, multi-dimensional time. *Int. J. Parallel Programming*, 21(6):389–420, December 1992.
- [9] Paul Feautrier. Towards automatic distribution. *Parallel Processing Letters*, 4(3):233–244, 1994.
- [10] Hans Michael Gerndt. *Automatic parallelization for distributed-memory multiprocessing systems*. PhD thesis, University of Bonn, 1989.
- [11] Erich Hecke. *Lectures on the Theory of Algebraic Numbers*. Springer-Verlag, 1981.
- [12] C.H. Huang and P. Sadayappan. Communication-free hyperplane partitioning of nested loops. In Banerjee, Gelernter, Nicolau, and Padua, editors, *Languages and Compilers for Parallel Computing*, volume 589 of *Lecture Notes in Computer Science*, pages 186–200. Springer Verlag, 1991.
- [13] Kathleen Knobe, Joan D. Lukas, and Guy L. Steele. Data optimization: Allocation of arrays to reduce communication on SIMD machines. *Journal of Parallel and Distributed Computing*, 8:102–118, 1990.
- [14] A. Korganoff and M. Pavel-Parvu. *Éléments de théorie des matrices carrées et rectangles en analyse numérique*. Dunod, 1966. (In French).
- [15] Jingke Li and Marina Chen. The data alignment phase in compiling programs for distributed memory machines. *Journal of Parallel and Distributed Computing*, 13:213–221, 1991.
- [16] Joan D. Lukas. Data locality for shared memory. In *Sixth SIAM Conf. Parallel Processing for Scientific Computing*, volume 2, pages 836–839. SIAM Press, 1993.
- [17] Joan D. Lukas and Kathleen Knobe. Data optimization and its effect on communication costs in MIMD fortran code. In Dongarra, Kennedy, Messina, Sorensen, and Voigt, editors, *Fifth SIAM Conf. Parallel Processing for Scientific Computing*, pages 478–483. SIAM Press, 1992.
- [18] Morris Newman. *Integral Matrices*. Academic Press, 1972.
- [19] Alexis Platonoff. *Contribution à la Distribution Automatique des Données pour Machines Massivement Parallèles*. PhD thesis, École Nationale Supérieure des Mines de Paris, March 1995.
- [20] W. Shang and Z. Shu. Data alignment of loop nests without nonlocal communications. In *Application Specific Array Processors*, pages 439–450. IEEE Computer Society Press, 1994.
- [21] Chau-Wen Tseng. *An Optimizing Fortran D Compiler for MIMD Distributed-Memory Machines*. PhD thesis, Rice University, 1993.
- [22] Michael Wolfe. The Tiny loop restructuring research tool. In H.D. Schwetman, editor, *International Conference on Parallel Processing*, volume II, pages 46–53. CRC Press, 1991.