

Memory Cost due to Anticipated Broadcast

Vincent Loechner and Catherine Mongenet

Laboratoire ICPS
Université Louis Pasteur de Strasbourg

October 15, 1999

1 Introduction

Over the last years many research works have focused on the parallelization of nested loops. To get an efficient parallelization, these techniques mainly focus on data alignment or on communication minimization [RS91, LC91, AL93, DR94, Fea94, BKK⁺95, DR95, CGO⁺96, Mon97, LM98]. The efficiency of a parallel solution not only depends on the communication cost, but also on the memory cost. This paper mainly focus on a symbolic evaluation of the memory cost.

The parallelization of a loop nest consists in determining a schedule function which states for any elementary computation the time when it is executed, and an allocation function which associates a virtual processor to each elementary computation. In order to be valid the parallelization of a loop nest relies on its dependences. Besides well studied uniform dependences, many scientific problems are expressed using loops with affine dependences. For instance, matrix computations such as triangulation use a diagonal element to compute all the elements in the same row or column. In a parallel implementation, this diagonal element is computed by a given processor (called *emission processor*), which has to send it to all the processors that use it (called *utilization processors*). This can be done using two communication schemes:

- pipeline: the data is sent from the emission processor to one of the utilization processors, and then propagated from neighbor to neighbor among the utilization processors ;
- broadcast: the data is simultaneously sent from the emission processor to all the utilization processors.

On a parallel architecture where broadcast is efficiently implemented, the second communication scheme should be chosen. Hence, it corresponds to a single communication primitive, whereas the first scheme requires as many communication primitive calls as utilization processors. This is in particular relevant when the network is a bus, as it occurs in a cluster of workstations connected by ethernet for example. On this type of architecture, since the cost of a broadcast

communication is equal to the cost of one point-to-point communication, the broadcast scheme, as opposed to the pipeline scheme, drastically reduces the network load. This approach is of course irrelevant on any parallel architecture where the broadcast is simulated by a pipeline, such as on systolic-like architectures. In the following, we assume that the broadcast is efficiently implemented and we focus on techniques to derive parallel loops using as much broadcast as possible.

The ideal broadcast communication scheme is the following: a given data is computed, it is then sent to all its utilization processors which use it right away. We call it *true broadcast*. For a given data, a true broadcast can only be implemented if all the computations that use it are scheduled at the same time step. Therefore true broadcast imposes strong conditions to be satisfied by the schedule function. For a problem characterized by several dependences, these conditions can not always be satisfied simultaneously. When they are not satisfied for a given dependence, it means that the utilization processors use the data at successive time steps. Even though pipeline could be used in this case, we choose to still use broadcast to reduce the communication cost. We call it *anticipated broadcast*. Since each utilization processor does not use the data immediately after its reception, it must store it until it finally uses it. This memory cost should be controlled in order to fit the hardware.

This work presents a method, based on a symbolic evaluation, to compute for given schedule and allocation functions, the memory cost related to each dependence in case of an anticipated broadcast. This symbolic evaluation is conducted in the framework of systems of parameterized affine recurrence equations which abstract single assignment loop nests. Using the dependence modeling in terms of utilization and emission sets [MCP94], classical affine per variable space-time transformations are computed [LM97]. In this approach, called the polytope model [Len93], the information can be represented by parameterized¹ convex polyhedra: utilization and emission sets, isotemporal spaces (i.e. polyhedra containing all the computation points scheduled at the same time step), etc. One of the techniques used to extract information from such polyhedra relies on counting the number of integer points it contains [CL98]. For instance, the potential parallelism at a given time step, i.e. the number of active processors at that time step, is the number of integer points in the corresponding isotemporal space. This symbolic evaluation is realized using Ehrhart polynomials [Cla96] which express this number as a function of its parameters.

The paper is organized as follows. It first summarizes the classical framework of systems of affine recurrence equations parallelization. It particularly recalls how to determine different space-time transformations. It then shows how to compute, for a given transformation, the memory cost using appropriate Ehrhart polynomials. It illustrates the approach on the Gaussian elimination problem using two different schedules. It finally gives some concluding remarks.

¹The parameters can be: the problem size parameters, a time step, a processor coordinates, an emission point, ...

2 Parallelization framework

A single assignment loop nest can be formalized by a system of parameterized affine recurrence equations (PARE). Such an equation is defined by:

$$X[z] = f(\dots, Y[g(z, p)], \dots) \quad z \in \mathcal{D}(p) \quad (Eq)$$

where X and Y are variable names. p is a vector of \mathbb{Z}^q defining the problem parameters. $\mathcal{D}(p) \subset \mathbb{Z}^n$ is a parameterized convex bounded polyhedron called the *iteration domain*, defined by a set of linear constraints. $g(z, p)$ is an affine function of z and p .

Dependence modeling

The dependence information is extracted from such a PARE in terms of utilization and emission sets, as presented in [MCP91, LM98].

The *utilization set* corresponding to occurrence $Y[z_0]$ in equation (Eq) is the set of points z where a computation $X[z]$ needs $Y[z_0]$. It is a convex polyhedron parameterized by z_0 and p :

$$\text{Util}_{Eq,Y}(z_0, p) = \{z \in \mathcal{D}(p) \mid g(z, p) = z_0\}.$$

Its vertices can be computed using the Loechner-Wilde algorithm [LW97], as function of z_0 and p . The geometric dimension of $\text{Util}_{Eq,Y}(z_0, p)$ is denoted by $\dim_{\text{Util}_{Eq,Y}}$. This set expresses that there is a dependence between point z_0 and each of the points of $\text{Util}_{Eq,Y}(z_0, p)$.

The set of all points z_0 which originate such a dependence is called the *emission set*:

$$\begin{aligned} \text{Emit}_{Eq,Y}(p) &= \{z_0 \in \mathbb{Z}^n \mid \exists z \in \mathcal{D}(p) \text{ such that } z_0 = g(z, p)\} \\ &= g(\mathcal{D}(p), p) \end{aligned}$$

In the following, we simplify notation $g(\mathcal{D}(p), p)$ in $g(\mathcal{D}(p))$.

Space-time transformation

In order to find a parallel solution to a set of PARE an affine space-time transformation has to be determined. Such a transformation is classically defined by two functions: an affine schedule function and an affine allocation function. In order to efficiently optimize the communications, we deal with *per variable* transformation. Such a transformation is expressed, for each variable X , by a full-column rank transformation matrix $T_X = \begin{pmatrix} \lambda_X \\ \sigma_X \end{pmatrix}$ and two constant vectors α_X and γ_X where:

- λ_X is a $d_\lambda \times n$ integer matrix defining the *multi-dimensional schedule*: $t_X(z) = \lambda_X z + \alpha_X$. The row vectors of λ_X are called the *schedule vectors*. They are orthogonal to a set of affine subspaces called *isotemporal spaces*. Each isotemporal space contains all the computation points of X executed at a given time step.

- σ_X is a $d_\sigma \times n$ integer matrix representing the *allocation* on a virtual processor space of dimension d_σ . The row vectors of σ_X generate the virtual processor space associated with X . The allocation function is defined by $alloc_X(z) = \sigma_X z + \gamma_X$.

Once a set of per variable transformations has been determined, the parallel code can be synthesized by applying these transformations to the PARE, as presented in [Ram92, CFR95] for a single transformation, and in [Alb98, QRW99] for per variable transformation.

For the transformations to be valid, matrices T_X must be full rank. Hence, it guarantees that no two computations are executed at the same time step by the same processor. Of course, not every full rank matrix is valid. Other conditions must or may be enforced. The first conditions that *must always* be satisfied are the *causal constraints*. These constraints state that a data has to be computed before its use. They are classically expressed for any dependence from Y to X of the form $X[z] \leftarrow Y[g(z, p)]$ by

$$t_X(z) \succ t_Y(z_0), \quad \forall z \in \text{Util}_{E_q, Y}(z_0, p)$$

where \succ is the lexicographical *greater than* operator.

Communication optimization

Other conditions may be added to optimize communications. As presented in earlier works [AL93, Fea94, DR95, Mon97, LM98], a *strong condition* can be computed so to map, for a given dependence, any utilization set on a single processor. It results into single communications from each emission processor to a unique utilization processor. This condition imposes constraints on the allocation matrices and on the schedule matrices. For many problems characterized by several dependences, this condition can not be satisfied for all of them. In such a case one approach is to apply the strong condition for as many dependences as possible, and then to optimize the residual communications [DR95]. Another approach is to combine the strong condition with other communication minimization conditions that do not impose such strong constraints on the transformation matrices. We propose two such conditions. The first one is a *weak condition* that enforces pipeline (also called *propagation*) of the data among the processors [LM97]. The second one takes advantage of broadcast communications, and in [LM98] various conditions on the schedule and the allocation functions are given to ensure local and/or normal broadcast. In the following, we focus on broadcast and we show how to communicate using either true or anticipated broadcast.

Symbolic evaluation

In order to evaluate the quality of a synthesized parallel solution, in particular its efficiency in term of latency, number of processors, communication cost and memory cost, symbolic evaluation must be conducted. The objects manipulated

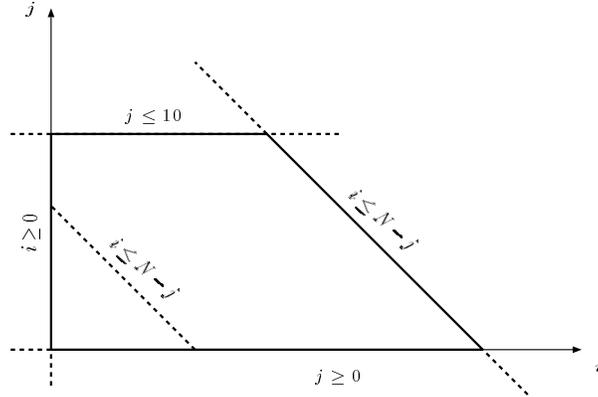


Figure 1: Parameterized polytope $\mathcal{D}(N) = \{(i, j) | 0 \leq i \leq N - j, 0 \leq j \leq 10\}$.

in the framework of PAREs are convex parameterized polyhedra. The parameters can be the problem parameters p , an emission point z_0 , a time step τ , a processor reference a , or any combination of them. The evaluation consists in counting the number of points in such a bounded polyhedron or polytope.

For example, the elementary computations of an iteration domain $\mathcal{D}(p)$, executed on a given processor a , are expressed by the following polytope:

$$\mathcal{E}(a, p) = \text{alloc}^{-1}(a) \cap \mathcal{D}(p)$$

where alloc^{-1} is the preimage of the allocation function. It defines the set of all the points whose image by function alloc is a . Therefore the intersection of this set with $\mathcal{D}(p)$ characterizes all the elementary computations executed on processor a . Counting the number of points of this polytope parameterized by a and p will give the number of computations executed on processor a , as a function of a and p .

Such a counting is realized using Ehrhart polynomials [Cla96]. Depending on the actual values of the parameters, the shape of the polytope may be different. Therefore an Ehrhart polynomial is defined on each *validity domain* of the parameters. For instance, consider the polytope defined by $\mathcal{D}(N) = \{(i, j) \in \mathbb{N}^2 | 0 \leq i \leq N - j, 0 \leq j \leq 10\}$. Depending on the value of parameter N (less or equal to 10 or greater than 10) this polyhedron has either 3 or 4 vertices, as shown in figure 1. The number of integer points contained in this polyhedron is defined by the following Ehrhart polynomial:

$$\begin{cases} \text{if } 0 \leq N \leq 10 & \frac{1}{2}N^2 + \frac{3}{2}N + 1 \\ \text{if } N \geq 10 & 11N - 44 \end{cases}$$

In this example each polynomial has rational coefficients, but the it is not always the case. In general, the coefficients are periodic numbers: their value are defined by a **case** function of the parameters *modulo* a constant. Such a polynomial is

called a pseudo-polynomial. For instance, the number of points in the polytope $\mathcal{D}(N) = \{k \in \mathbb{Z} | 0 \leq k \leq N/2\}$ is:

$$EP(N) = N + [1, 0]_N.$$

where $[1, 0]_N$ is a periodic number defined by:

$$\begin{cases} 1 & \text{if } N \bmod 2 = 0 \\ 0 & \text{if } N \bmod 2 = 1 \end{cases}$$

The computation of Ehrhart polynomials has been implemented in the *PolyLib* library (<http://icps.u-strasbg.fr/PolyLib>) and the results presented in section 5 are computed in this library.

3 Broadcast communication schemes

The ideal broadcast communication scheme occurs when all the utilization processors use a given data at the same time step. We call it a *true broadcast*. This imposes strong constraints on the schedule. In the simple case of a broadcast along a one-dimensional utilization set, all its points must be scheduled at the same time step [MCP91]. In the general case of multi-dimensional utilization sets, the condition on the schedule function is given in [LM98].

The memory cost associated with true broadcast communication is minimal. For instance, in the case of a broadcast along a one-dimensional utilization set, each utilization processor needs only one memory location to hold the data.

However, due to the strong conditions it imposes on the schedule functions, in many cases true broadcast can not be enforced for all the dependences. In this case, for those dependences where true broadcast can not be applied, we propose to use *anticipated broadcast*. This means that even though broadcast is not necessary (relatively to the schedule) we still broadcast the data. We do so in order to have only one communication as opposed to a set of neighbour to neighbor communications if pipeline would be used. Notice that this type of broadcast can always be applied, without any constraints on the schedule.

For a given utilization processor, there is a delay between the time step it receives the data and the time step it first uses it. Therefore, the data must be stored from the time step it is received until the one corresponding to its last use. We call this interval the *lifetime* of the data. During this interval, the processor may receive other data to be used later. It therefore needs memory locations for all of these data. Of course, in order to optimize memory usage, one must not retain a memory location for a data after its last use. This memory location should be freed and allocated to another data used by the processor at that time step [WR97].

For a problem with several dependences, there is an important degree of freedom in the choices of space-time transformations and communication schemes. In order to determine a large spectrum of efficient solutions one should apply a heuristic that proposes various space-time transformations and determines for

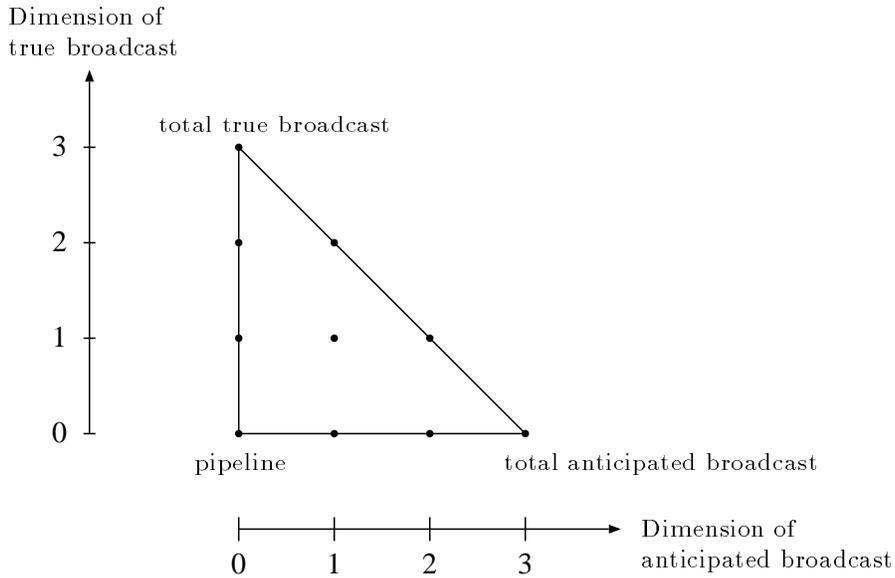


Figure 2: Communications implementation: dimension of anticipated broadcast.

each of them a communication cost and a memory cost, as presented in [LM98]. Such a heuristic analyses successively each of the dependences of the problem, and determines the constraints on the schedule and the allocation functions relatively to these various communication schemes.

The degree of freedom for a given dependence is pictured in figure 2 for a 3-dimensional utilization set. One end of the spectrum is characterized by a total true broadcast (3-dimensional true broadcast). This obviously corresponds to a strongly constrained schedule function: all the points of this 3-dimensional utilization set are scheduled at the same time step. This is the upper left point on the figure. If we relax the constraints on the schedule, the dimension of true broadcast decreases. A 2-dimensional true broadcast is represented by the points of ordinate 2. For the third dimension of the utilization set where there is no true broadcast, one may either choose to use pipeline (it corresponds to the left-hand side point), or to use anticipated broadcast (right-hand side point). The points of ordinate 0 correspond to space-time transformations avoiding true broadcast. In such a case there is still a large degree of freedom, from pipeline in all three dimensions (left-hand side point) to total (3-dimensional) anticipated broadcast (right-hand side point). The intermediate points correspond to partial anticipated broadcast communications associated with pipeline on the remaining dimensions.

In [LM98] the memory cost for anticipated broadcast was evaluated as a naive upper bound: the lifetime intervals were not taken into account, but

considered as infinite. Our objective here is to symbolically compute the exact memory cost for a complete anticipated broadcast (right-hand side points of figure 2), as a function of the processor coordinates and the problem parameters.

4 Memory cost evaluation for anticipated broadcast

An anticipated broadcast can be implemented using two symmetric memory management schemes:

1. the data is sent by the emission processor immediately after its computation to all its utilization processors². In this case the memory cost is entirely carried by the utilization processors.
2. the data is stored on the emission processor until its first use by one of the utilization processors. At this time step it is simultaneously sent to all the utilization processors.

The memory required on a processor to store the data for a given dependence $X \leftarrow Y$ is obviously strongly related to the allocation and the schedule functions associated with both variables. When any processor receives a single item of variable Y , i.e. there is a single utilization set mapped on any processor, the memory cost is one. The interesting situation occurs when a given utilization processor receives successively different items of variable Y (i.e. when several utilization sets are mapped on this processor).

A given item $Y[z_0]$ must be stored on utilization processor a between the following time steps:

- $\tau_{min}(z_0, p, a)$: the time step when the anticipated broadcast is realized.
 - In the first memory management scheme, the broadcast is realized immediately after the data is computed. In this case $\tau_{min}(z_0, p, a) = t_Y(z_0) + \vec{1}$ where $\vec{1}$ is the smallest time unit, i.e. d_λ -vector $(0, \dots, 0, 1)$ for a d_λ -dimensional schedule.
 - In the second memory management scheme, the broadcast is realized immediately before the first use of the data. Since the schedule is affine, the first use of the data occurs at a vertex of the utilization set. Let $v_i(z_0, p)$ denote the vertices of $\text{Util}_Y(z_0, p)$. It follows that $\tau_{min}(z_0, p, a) = \min_i(t_X(v_i(z_0, p, a)))$.
- $\tau_{max}(z_0, p, a)$: the time step corresponding to the last use on this utilization processor. Let $\text{UP}(z_0, p, a) = \text{alloc}_Y^{-1}(a) \cap \text{Util}_Y(z_0, p)$ be the set of utilization points mapped on processor a . It is a convex parameterized

²Whatever the dimension of the utilization set, the broadcast will be realized over all the utilization processors. An alternative would be to cut this global broadcast into successive broadcasts of smaller dimension.

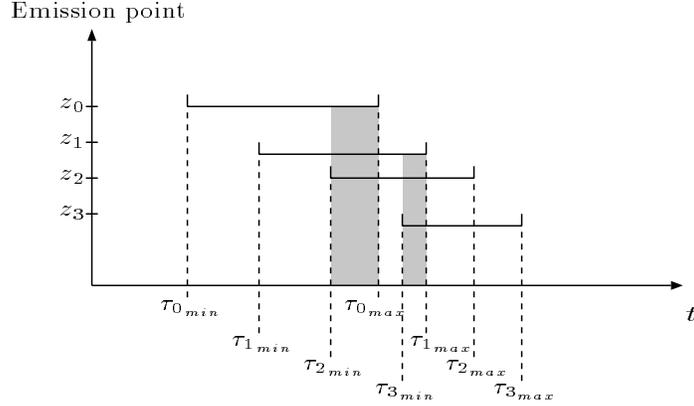


Figure 3: Lifetime intervals for a given processor.

polyhedron characterized by a set of parameterized vertices $v'_i(z_0, p, a)$. It follows that $\tau_{max}(z_0, p, a) = \max_i(t_X(v'_i(z_0, p, a)))$.

We call *lifetime* the time interval between $\tau_{min}(z_0, p, a)$ and $\tau_{max}(z_0, p, a)$. Notice that when the schedule function is multidimensional, we consider the lexicographic order on the time step vectors. Therefore, the lifetime is a union of parameterized polyhedra.

Definition 1 For a given dependence from Y to X , the lifetime of $Y[z_0]$ on processor a is:

$$\mathcal{L}(z_0, p, a) = \{\tau | \tau_{min}(z_0, p, a) \leq \tau \leq \tau_{max}(z_0, p, a)\}.$$

When a given utilization processor receives successive data, they must be stored during their respective lifetimes. Two cases may occur:

- the various lifetimes are all disjoint, in which case the memory cost is one;
- the lifetimes do intersect, in which case the memory cost is the maximum at any time step of the number of intersecting sets. Figure 3 represents four lifetime intervals related to four emission points for a given processor. The maximal number of intersecting sets is 3, between $\tau_{2_{min}}$ and $\tau_{0_{max}}$ or between $\tau_{3_{min}}$ and $\tau_{1_{max}}$. Thus the maximal number of simultaneously alive data, that is the memory cost, is 3.

In the general case (when the lifetimes intersect) the amount of memory required on a processor is computed according to the following theorem.

Theorem 1 For a given dependence from Y to X , the memory required on a utilization processor a is:

$$mem(p, a) = \underset{z_0 \in \text{Emit}_Y(p)}{Max} (EP(g(t_X^{-1}(\mathcal{L}(z_0, p, a)) \cap \mathcal{D}_{Eq}(p) \cap alloc_X^{-1}(a))).$$

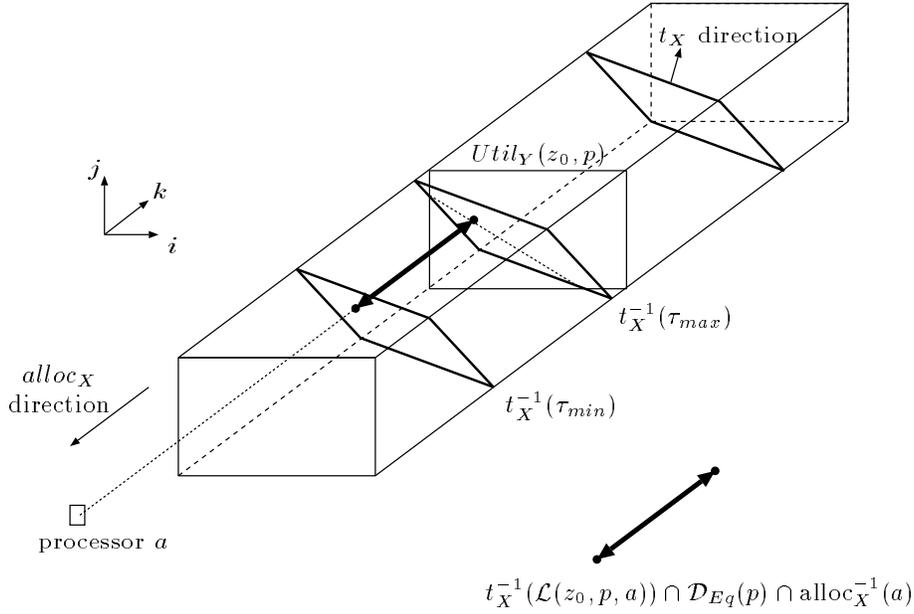


Figure 4: Active utilization points during the lifetime of data $Y[z_0]$ on processor a .

Proof.

For a given point $z_0 \in \text{Emit}_Y(p)$, the lifetime of $Y[z_0]$ on processor a is $\mathcal{L}(z_0, p, a)$. This set represents the time steps at which $Y[z_0]$ must be stored on utilization processor a . The inverse image by t_X intersected with the iteration domain, $t_X^{-1}(\mathcal{L}(z_0, p, a)) \cap \mathcal{D}_{Eq}(p)$, is the set of all the points of the domain that are active during this lifetime. The intersection of this set with $\text{alloc}_X^{-1}(a)$ restricts it to the points mapped on processor a .

Figure 4 shows such a set of points. The main parallelepiped represents the domain $\mathcal{D}_{Eq}(p)$. For a given emission point³ z_0 , its corresponding utilization set $\text{Util}_Y(z_0, p)$ is shown. The isotemporal spaces $t_X^{-1}(\tau_{min})$ and $t_X^{-1}(\tau_{max})$ are represented⁴. Notice that τ_{min} depends on the memory management scheme as described above. The bold segment on the figure represents $\mathcal{S}(z_0, p, a) = t_X^{-1}(\mathcal{L}(z_0, p, a)) \cap \mathcal{D}_{Eq}(p) \cap \text{alloc}_X^{-1}(a)$. This segment, called the *utilization segment*, contains the utilization points whose corresponding data must be stored on processor a at the $Y[z_0]$ transmission time step.

³Notice that this point belongs to the definition domain of variable Y , therefore it does not appear on the picture which focuses on variable X .

⁴For the sake of clarity, we represent only the part of the isotemporal spaces intersecting the domain.

In this example, there is only one point per utilization set belonging to the utilization segment. However, in the general case, there could be several of them in the same utilization set. Since they correspond to only one data, they require only one memory location. In order to compute the amount of memory required when $Y[z_0]$ is sent, we must count the number of utilization sets, or equivalently the number of corresponding emission points. To get these emission points we apply function g to the utilization segment. We compute the Ehrhart polynomial of the resulting set to get the memory cost related to z_0 :

$$EP(g(\mathcal{S}(z_0, p, a)))$$

In order to get the final memory cost on processor a , we compute the maximum over $z_0 \in \text{Emit}_Y(p)$ of this expression.

□

In the simple case where the lifetimes do not intersect, theorem 1 simplifies as $mem(p, a) = 1$. Hence, each utilization segment is contained in a single utilization set $\text{Util}_Y(z_0, p)$. Therefore its image by function g is the single emission point z_0 , and $EP(g(\mathcal{S}(z_0, p, a))) = 1$ for any $z_0 \in \text{Emit}_Y(p)$.

If the first memory management scheme is used, i.e. if the emission processor sends the data to the utilization processors immediately after its computation, then theorem 1 gives the total amount of memory required on any processor for this dependence. If the second memory management scheme is used, i.e. if the data is stored on the emission processor until its first use by one of the utilization processors, then one must compute the amount of memory required on an emission processor. This is expressed by the following theorem.

Theorem 2 *For a given dependence from Y to X the memory required on an emission processor a , when using the second memory management scheme, is:*

$$mem(p, a) = \underset{z_0 \in \text{Emit}_Y(p)}{\text{Max}} (EP(g(t_X^{-1}(\mathcal{L}'(z_0, p, a)) \cap \mathcal{D}_{E_q}(p) \cap \text{alloc}_X^{-1}(a))))$$

where $\mathcal{L}'(z_0, p, a) = \{\tau | t_Y(z_0) < \tau \leq \tau_{min}(z_0, p, a)\}$ is the time interval during which the data is stored on the emission processor.

The proof is similar to the one of theorem 1 by substituting set \mathcal{L} by \mathcal{L}' .

These two theorems characterize, for a given dependence, the memory required on both the emission and the utilization processors. For a problem, they have to be applied to all its dependences. This is illustrated in the next section on an example.

5 Example: Gaussian elimination

Let us consider the gaussian elimination problem, which triangulates an $N \times (N + 1)$ matrix m . The resulting upper-triangular matrix is denoted by M . It

Dependence		Util(i_0, j_0, k_0, N)
(2)	$A[i, k, k - 1]$	$\left\{ \binom{i_0}{j} \middle k_0 + 1 \leq j \leq N + 1 \right\}$
(3)	$A[k, j, k - 1]$	$\left\{ \binom{i}{j_0} \middle k_0 + 1 \leq i \leq N \right\}$
(4)	$A[k, k, k - 1]$	$\left\{ \binom{i}{j} \middle k_0 + 1 \leq i \leq N, k_0 + 1 \leq j \leq N + 1 \right\}$

Table 1: Utilization sets.

is defined by the following system of affine recurrence equations:

$$A[i, j, 0] = m_{i,j} \quad (1)$$

$$\mathcal{D}_1(N) = \{(i, j) \in \mathbb{Z}^2 \mid 1 \leq i \leq N; 1 \leq j \leq N + 1\}$$

$$A[i, j, k] = A[i, j, k - 1] - A[i, k, k - 1] * A[k, j, k - 1] / A[k, k, k - 1] \quad (2)$$

$$\mathcal{D}_2(N) = \{(i, j, k) \in \mathbb{Z}^3 \mid k + 1 \leq i \leq N; k + 1 \leq j \leq N + 1; 1 \leq k \leq N - 1\}$$

$$M[i, j] = A[i, j, i - 1] \quad (3)$$

$$\mathcal{D}_3(N) = \{(i, j) \in \mathbb{Z}^2 \mid 1 \leq i \leq N; i \leq j \leq N + 1\}$$

The definition domain of variable A is:

$$\mathcal{D}_A(N) = \{(i, j, k) \in \mathbb{Z}^3 \mid k + 1 \leq i \leq N; k + 1 \leq j \leq N + 1; 0 \leq k \leq N - 1\}$$

In the following, we focus on computation equation (2) defined over iteration domain $\mathcal{D}_2(N)$. The first dependence is uniform and is therefore of no interest in the discussion. It always requires one memory location which should be added to the results presented hereunder, where we only consider the last three dependence, numbered (2) to (4).

In order to find a schedule and an allocation function a dependence analysis has first to be conducted. The utilization and emission sets related to each dependence are computed using the OPERA environment [LM96], and summarized in tables 1 and 2. Using the communication minimization algorithm presented in [LM98] one can determine efficient space-time transformations. When restricting to one-dimensional schedules, the schedule matrix λ_A can be chosen in the following interval:

- $\lambda_A = \begin{pmatrix} 0 & 0 & 1 \end{pmatrix}$: this is the fastest affine schedule relatively to the dependences. It enforces true broadcast for all the dependences. When applying the results presented in this paper, one of course obtains a total amount of memory equal to 3.

Dependence		Emit(N)
(2)	$A[i, k, k - 1]$	$\left\{ \binom{i}{k+1} \mid k+2 \leq i \leq N, 0 \leq k \leq N-2 \right\}$
(3)	$A[k, j, k - 1]$	$\left\{ \binom{k+1}{j} \mid k+2 \leq j \leq N+1, 0 \leq k \leq N-2 \right\}$
(4)	$A[k, k, k - 1]$	$\left\{ \binom{k+1}{k+1} \mid 0 \leq k \leq N-2 \right\}$

Table 2: Emission sets.

- $\lambda_A = (1 \ 1 \ 1)$: this schedule is the fastest affine schedule without any true broadcast.

In the following, we focus on the first memory management scheme described above, in which the broadcast is realized immediately after the data is computed. We consider two schedule matrices involving anticipated broadcast: $(1 \ 1 \ 1)$ where all three dependences may be realized using anticipated broadcast; $(0 \ 1 \ 1)$ where true broadcast is partially required, and anticipated broadcast can be used. The memory cost in this case is obviously lower than in the first case.

In both cases, we choose the same allocation function defined by matrix $\sigma_A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}$. We denote by $\begin{pmatrix} a_0 \\ a_1 \end{pmatrix}$ the coordinates of a given processor. In order to compute the lifetime, the sets UP of utilization points mapped on a given processor a are shown in table 3. Notice that each set $UP(i_0, j_0, k_0, N, a_0, a_1)$ is reduced to a single point which is therefore its unique vertex, denoted by $v'_1(i_0, j_0, k_0, N, a_0, a_1)$.

5.1 First schedule: $\lambda_A = (1 \ 1 \ 1)$

Using the first memory management scheme, the extremal values τ_{min} and τ_{max} are, for all three dependences:

$$\begin{aligned} \tau_{min}(i_0, j_0, k_0, N, a_0, a_1) &= t_A(i_0, j_0, k_0) + 1 = i_0 + j_0 + k_0 + 1 \\ \tau_{max}(i_0, j_0, k_0, N, a_0, a_1) &= t_A(v'_1(i_0, j_0, k_0, N, a_0, a_1)) = a_0 + a_1 + k_0 + 1 \end{aligned}$$

Dependence		UP $(i_0, j_0, k_0, N, a_0, a_1)$
(2)	$A[i, k, k - 1]$	$\left\{ \binom{a_0}{a_1}{k_0 + 1} \mid i_0 = a_0, k_0 + 1 \leq a_1 \leq N + 1 \right\}$
(3)	$A[k, j, k - 1]$	$\left\{ \binom{a_0}{a_1}{k_0 + 1} \mid a_1 = j_0, k_0 + 1 \leq a_0 \leq N \right\}$
(4)	$A[k, k, k - 1]$	$\left\{ \binom{a_0}{a_1}{k_0 + 1} \mid k_0 + 1 \leq a_0 \leq N, k_0 + 1 \leq a_1 \leq N + 1 \right\}$

Table 3: Utilization points mapped on processor a .

The lifetime of $A[i_0, j_0, k_0]$ on processor $\binom{a_0}{a_1}$ for the three dependences⁵ is:

$$\mathcal{L}(i_0, j_0, k_0, N, a_0, a_1) = \{\tau \mid i_0 + j_0 + k_0 + 1 \leq \tau \leq a_0 + a_1 + k_0 + 1\}.$$

The amount of memory required on any utilization processor is given by theorem 1. The utilization segment used in the theorem is defined by:

$$\begin{aligned} \mathcal{S}(i_0, j_0, k_0, N, a_0, a_1) &= t_A^{-1}(\mathcal{L}(i_0, j_0, k_0, N, a_0, a_1)) \cap \mathcal{D}_2(N) \cap \text{alloc}_A^{-1}(a_0, a_1) \\ &= \left\{ \binom{i}{j}{k} \mid \begin{array}{l} i_0 + j_0 + k_0 + 1 \leq i + j + k \leq a_0 + a_1 + k_0 + 1, k + 1 \leq i, \\ i \leq N, k + 1 \leq j \leq N + 1, 1 \leq k \leq N - 1, i = a_0, j = a_1 \end{array} \right\} \\ &= \left\{ \binom{a_0}{a_1}{k} \mid i_0 + j_0 + k_0 + 1 \leq a_0 + a_1 + k, 1 \leq k \leq k_0 + 1 \right\}. \end{aligned}$$

Since the memory cost is computed for a given dependence, the expression of the utilization segment depends on the validity domains of the parameters for each dependence, which are:

- for dependence (2): $k_0 + 2 \leq i_0 \leq N, j_0 = k_0 + 1, 0 \leq k_0 \leq N - 2, a_0 = i_0, k_0 + 1 \leq a_1 \leq N + 1, 1 \leq N$.
- for dependence (3): $i_0 = k_0 + 1, k_0 + 2 \leq j_0 \leq N + 1, 0 \leq k_0 \leq N - 2, k_0 + 1 \leq a_0 \leq N, a_1 = j_0, 1 \leq N$.
- for dependence (4): $i_0 = k_0 + 1, j_0 = k_0 + 1, 0 \leq k_0 \leq N - 2, k_0 + 1 \leq a_0 \leq N, k_0 + 1 \leq a_1 \leq N + 1, 1 \leq N$.

Let us now examine successively each of the three dependences.

⁵The formulation is identical for all three dependences, however the constraints on the parameters $i_0, j_0, k_0, N, a_0, a_1$ vary from one dependence to another.

5.1.1 Dependence (2)

$$g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1)) = \left\{ \left(\begin{array}{c} a_0 \\ k \\ k-1 \end{array} \right) \middle| \begin{array}{l} i_0 + j_0 + k_0 + 1 \leq a_0 + a_1 + k, \\ k \leq k_0 + 1 \end{array} \right\}$$

$$EP(g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1))) = \begin{cases} k_0 + 1 & \text{if } 0 \leq k_0 \leq \frac{a_1 - 1}{2} \\ a_1 - k_0 & \text{if } \frac{a_1 - 1}{2} \leq k_0 \leq a_1 - 2 \end{cases}$$

$$\begin{aligned} mem(N, a_0, a_1) &= \underset{i_0, j_0, k_0 \in \text{Emit}_Y(N)}{Max} (EP(g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1)))) \\ &= \left\lfloor \frac{a_1 + 1}{2} \right\rfloor \end{aligned}$$

5.1.2 Dependence (3)

$$g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1)) = \left\{ \left(\begin{array}{c} k \\ a_1 \\ k-1 \end{array} \right) \middle| \begin{array}{l} i_0 + j_0 + k_0 + 1 \leq a_0 + a_1 + k, \\ k \leq k_0 + 1 \end{array} \right\}$$

$$EP(g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1))) = \begin{cases} k_0 + 1 & \text{if } k_0 \leq \frac{a_0 - 1}{2} \\ a_0 - k_0 & \text{if } k_0 \geq \frac{a_0 - 1}{2} \end{cases}$$

$$\begin{aligned} mem(N, a_0, a_1) &= \underset{i_0, j_0, k_0 \in \text{Emit}_Y(N)}{Max} (EP(g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1)))) \\ &= \left\lfloor \frac{a_0 + 1}{2} \right\rfloor \end{aligned}$$

5.1.3 Dependence (4)

$$g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1)) = \left\{ \left(\begin{array}{c} k \\ k \\ k-1 \end{array} \right) \middle| \begin{array}{l} i_0 + j_0 + k_0 + 1 \leq a_0 + a_1 + k, \\ k \leq k_0 + 1 \end{array} \right\}$$

$$EP(g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1))) = \begin{cases} k_0 + 1 & \text{if } k_0 \leq \frac{a_0 + a_1 - 2}{3} \\ a_0 + a_1 - 2k_0 - 1 & \text{if } k_0 \geq \frac{a_0 + a_1 - 2}{3} \end{cases}$$

$$\begin{aligned} mem(N, a_0, a_1) &= \underset{i_0, j_0, k_0 \in \text{Emit}_Y(N)}{Max} (EP(g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1)))) \\ &= \left\lfloor \frac{a_0 + a_1 + 1}{3} \right\rfloor \end{aligned}$$

5.1.4 Memory required on any processor

The total amount of memory on a given processor is the sum on all three dependence of $mem(N, a_0, a_1)$:

$$mem(N, a_0, a_1) = \left\lfloor \frac{a_1 + 1}{2} \right\rfloor + \left\lfloor \frac{a_0 + 1}{2} \right\rfloor + \left\lfloor \frac{a_0 + a_1 + 1}{3} \right\rfloor$$

Applying the allocation function to the iteration domain $\mathcal{D}_2(N)$, the active virtual processor space is:

$$alloc(\mathcal{D}_2(N)) = \{(a_0, a_1) \in \mathbb{Z}^2 | 2 \leq a_0 \leq N, 2 \leq a_1 \leq N + 1\}.$$

According to the above formula the smallest amount of memory is 3 and it is obtained for processor of coordinates (2, 2). The largest amount of memory is upper bounded by $\frac{5}{3}N + 2$ for processor $(N, N + 1)$. This is the total amount of memory required when all three dependences are realized as total anticipated broadcast. If this memory cost is too expensive one can of course realize some of these communications as propagation.

5.2 Second schedule: $\lambda_A = (0 \ 1 \ 1)$

Using the first memory management scheme, the extremal values τ_{min} and τ_{max} are, for all three dependences:

$$\tau_{min}(i_0, j_0, k_0, N, a_0, a_1) = t_A(i_0, j_0, k_0) + 1 = j_0 + k_0 + 1 \quad (4)$$

$$\tau_{max}(i_0, j_0, k_0, N, a_0, a_1) = t_A(v'_1(i_0, j_0, k_0, N, a_0, a_1)) = a_1 + k_0 + 1 \quad (5)$$

The lifetime of $A[i_0, j_0, k_0]$ on processor $\begin{pmatrix} a_0 \\ a_1 \end{pmatrix}$ for the three dependences is:

$$\mathcal{L}(i_0, j_0, k_0, N, a_0, a_1) = \{\tau | j_0 + k_0 + 1 \leq \tau \leq a_1 + k_0 + 1\}$$

The utilization segment used in the theorem is defined by:

$$\begin{aligned} \mathcal{S}(i_0, j_0, k_0, N, a_0, a_1) &= t_A^{-1}(\mathcal{L}(i_0, j_0, k_0, N, a_0, a_1)) \cap \mathcal{D}_2(N) \cap alloc_A^{-1}(a_0, a_1) \\ &= \left\{ \begin{pmatrix} i \\ j \\ k \end{pmatrix} \middle| \begin{array}{l} j_0 + k_0 + 1 \leq j + k \leq a_1 + k_0 + 1, k + 1 \leq i \leq N, \\ k + 1 \leq j \leq N + 1, 1 \leq k \leq N - 1, i = a_0, j = a_1 \end{array} \right\} \\ &= \left\{ \begin{pmatrix} a_0 \\ a_1 \\ k \end{pmatrix} \middle| j_0 + k_0 + 1 \leq a_1 + k, k \leq k_0 + 1, 1 \leq k \leq N - 1 \right\} \end{aligned}$$

The validity domains of the parameters for each dependence are the same as in subsection 5.1 since the allocation function is the same. Let us now examine successively each of the three dependences.

5.2.1 Dependence (2)

$$\begin{aligned}
g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1)) &= \left\{ \binom{a_0}{k} \middle| j_0 + k_0 + 1 \leq a_1 + k, k \leq k_0 + 1 \right\} \\
EP(g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1))) &= \begin{cases} k_0 + 1 & \text{if } 0 \leq k_0 \leq \frac{a_1 - 1}{2} \\ a_1 - k_0 & \text{if } \frac{a_1 - 1}{2} \leq k_0 \leq \min(a_0, a_1) - 2 \end{cases} \\
mem(N, a_0, a_1) &= \underset{i_0, j_0, k_0 \in \text{Emit}_Y(N)}{Max} (EP(g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1)))) \\
&= \left\lfloor \frac{a_1 + 1}{2} \right\rfloor
\end{aligned}$$

5.2.2 Dependence (3)

$$\begin{aligned}
g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1)) &= \left\{ \binom{k}{a_1} \middle| j_0 + k_0 + 1 \leq a_1 + k, k \leq k_0 + 1 \right\} \\
EP(g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1))) &= 1
\end{aligned}$$

This result is not surprising since this dependence is implemented as a true broadcast due to the schedule. Therefore the memory required is reduced to 1 location.

$$mem(N, a_0, a_1) = 1$$

5.2.3 Dependence (4)

$$\begin{aligned}
g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1)) &= \left\{ \binom{k}{k} \middle| j_0 + k_0 + 1 \leq a_1 + k, k \leq k_0 + 1 \right\} \\
EP(g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1))) &= \begin{cases} k_0 + 1 & \text{if } 0 \leq k_0 \leq \frac{a_1 - 1}{2} \\ a_1 - k_0 & \text{if } \frac{a_1 - 1}{2} \leq k_0 \leq \min(a_0, a_1) - 2 \end{cases} \\
mem(N, a_0, a_1) &= \underset{i_0, j_0, k_0 \in \text{Emit}_Y(N)}{Max} (EP(g(\mathcal{S}(i_0, j_0, k_0, N, a_0, a_1)))) \\
&= \left\lfloor \frac{a_1 + 1}{2} \right\rfloor
\end{aligned}$$

5.2.4 Memory required on any processor

The total amount of memory on a given processor is the sum on all three dependence of $mem(N, a_0, a_1)$:

$$\begin{aligned}
mem(N, a_0, a_1) &= \left\lfloor \frac{a_1 + 1}{2} \right\rfloor + 1 + \left\lfloor \frac{a_1 + 1}{2} \right\rfloor \\
&= 2 \left\lfloor \frac{a_1 + 1}{2} \right\rfloor + 1
\end{aligned}$$

Applying the allocation function to the iteration domain $\mathcal{D}_2(N)$, the active virtual processor space is:

$$alloc(\mathcal{D}_2(N)) = \{(a_0, a_1) \in \mathbb{Z}^2 \mid 2 \leq a_0 \leq N, 2 \leq a_1 \leq N + 1\}.$$

According to the above formula the smallest amount of memory is 3 and it is obtained for processor of coordinates (2, 2). The largest amount of memory is upper bounded by $N + 2$ for the processors of ordinate $N + 1$. As compared to the result obtained with schedule (1 1 1), this result is better (N versus $\frac{5}{3}N$). This is mainly due to dependence (3), which is implemented by a true broadcast.

6 Conclusion

The work presented in this paper evaluates the memory cost related to dependences whose related communications are implemented using anticipated broadcast. For a given set of per variable space-time transformations, the exact amount of memory required on any processor is computed as a function of the program parameters and of the processor coordinates. This computation is realized using Ehrhart polynomials, which give the number of integer points contained in a parameterized convex polytope. The amount of memory we compute for a given processor is optimal: it corresponds to the minimum number of memory locations required and therefore relies on memory reuse. Hence, as soon as a data is no longer used by a processor, its memory location is released and allocated to another data. In [WR97] Wilde and Rajopadhye also deal with memory reuse based on lifetime information, but do not take advantage of the usage information (called usage table in their paper) to integrate broadcast communications. They define a memory allocation function which optimizes memory reuse, but do not give any method to compute the actual memory cost.

In this paper we only deal with total anticipated broadcast, that is for a multi-dimensional utilization set, we only use broadcast communications, either true or anticipated. An alternative could be to use, for dimensions where true broadcast does not apply, anticipated broadcast on some of the dimensions and propagation on the remaining ones. A first extension of this work will be to deal with partial anticipated broadcast: use simultaneously anticipated broadcast (in some of the dimensions) and propagation (in the other ones).

This work is an element of a global algorithm for efficient loop nests parallelization. The whole process consists in finding an optimal space-time transformation, according to various optimization criteria such as latency, communication cost, or memory cost. Latency optimization criteria are well known, as presented by Feautrier in [Fea92a, Fea92b]. An algorithm dealing with communication optimization has been presented in [LM97, LM98]. We propose here an approach to evaluate the memory cost related to anticipated broadcast. These results should of course be integrated in the whole process. Notice that, at that stage, memory cost is evaluated for a given space-time transformation. A further extension of this work will consist in finding appropriate conditions to determine

memory-optimal space-time transformations, as it is done in [LM97, LM98] for communication optimization.

For a given problem, the various optimization criteria may be antinomial: using anticipated broadcast reduces the communication cost⁶ but increases memory cost. Therefore, a trade-off between these optimizations must be found in order to find globally efficient transformations. The approach developed in [LM98] should be extended to take advantage of anticipated broadcast. The cost function should then integrate latency and required amount of memory, besides the number of communications.

References

- [AL93] J.M. Anderson and M.S. Lam. Global optimizations for parallelism and locality on scalable parallel machines. *ACM Sigplan Notices*, 28(6):112–125, 1993.
- [Alb98] O. Albiez. Parcours d’une superposition de domaines. In D. Mery and G.R. Perrin, editors, *RenPar’10*, 1998.
- [BKK⁺95] D. Bau, I. Kodukula, V. Kotlyar, K. Pingali, and P. Stodghill. Solving alignment using elementary linear algebra. In *Languages and Compilers for Parallel Computing*. LNCS 892, Springer-Verlag, 1995.
- [CFR95] J.-F. Collard, P. Feautrier, and T. Risset. Construction of DO loops from systems of affine constraints. *Parallel Processing Letters*, 5(3):421–436, 1995.
- [CGO⁺96] S. Chatterjee, J.R. Gilbert, L. Oliker, R. Schreiber, and T.J. Sheffler. Algorithms for automatic alignment of arrays. *Journal of Parallel and Distributed Computing*, 38(2):145–157, November 1996.
- [CL98] Ph. Clauss and V. Loechner. Parametric analysis of polyhedral iteration spaces. *Journal of VLSI Signal Processing*, 19(2):179–194, July 1998.
- [Cla96] Ph. Clauss. Counting solutions to linear and nonlinear constraints through Ehrhart polynomials: Applications to analyze and transform scientific programs. In *10th ACM Int. Conf. on Supercomputing, Philadelphia*, 1996.
- [DR94] A. Darté and Y. Robert. On the alignment problem. *Parallel Processing Letters*, 4(3):259–270, 1994.
- [DR95] M. Dion and Y. Robert. Mapping affine loop nests : new results. In B. Hertzberger and G. Serazzi, editors, *Int. Conf. on High-Performance Computing and Networking*, pages 184–189. Springer Verlag, May 1995. Milan, Italy.

⁶on appropriate architectures where broadcast is efficiently implemented.

- [Fea92a] P. Feautrier. Some efficient solutions to the affine scheduling problem, part 1 : one dimensional time. *Int. J. of Parallel Programming*, 21(5):313–348, October 1992.
- [Fea92b] P. Feautrier. Some efficient solutions to the affine scheduling problem, part 2 : multidimensional time. *Int. J. of Parallel Programming*, 21(6), December 1992.
- [Fea94] P. Feautrier. Towards automatic distribution. *Parallel Processing Letters*, 4(3):233–244, 1994.
- [LC91] J. Li and M. Chen. The data alignment phase in compiling programs for distributed-memory machines. *Journal of Parallel and Distributed Computing*, 13(2), 1991.
- [Len93] C. Lengauer. Loop parallelization in the polytope model. In *CONCUR 93*, 1993.
- [LM96] V. Loechner and C. Mongenet. Opera : A toolbox for loop parallelization. In Jelly, Gordon, and Croll, editors, *First International Workshop on Software Engineering for Parallel and Distributed Systems*, pages 134–145. Chapman&Hall, March 1996.
- [LM97] V. Loechner and C. Mongenet. Solutions to the communication minimization problem for affine recurrence equations. In Lengauer, Griehl, and Gortlatch, editors, *Europar'97*, pages 328–337. Springer-Verlag, August 1997. LNCS 1300.
- [LM98] V. Loechner and C. Mongenet. Communication optimization for affine recurrence equations using broadcast and locality. Technical Report RR 98-03, ICPS, <http://icps.u-strasbg.fr/pub-98/>, 1998. To appear in *Int. J. of Parallel Programming*, February 2000.
- [LW97] V. Loechner and D. K. Wilde. Parameterized polyhedra and their vertices. *International Journal of Parallel Programming*, 25(6), December 1997.
- [MCP91] C. Mongenet, Ph. Clauss, and G.R. Perrin. A geometrical coding to compile affine recurrence equations on regular arrays. In *5th International Parallel Processing Symposium, IPPS'91*, 1991.
- [MCP94] C. Mongenet, Ph. Clauss, and G.R. Perrin. Geometrical tools to map systems of affine recurrence equations on regular arrays. *Acta Informatica*, 31:137–160, 1994.
- [Mon97] C. Mongenet. Affine dependence classification for communications minimization. *International Journal of Parallel Programming*, 25(6), December 1997.

- [QRW99] F. Quilleré, S. Rajopadhye, and D. Wilde. Generation of efficient nested loops from polyhedra. Research report, IRISA, Rennes, France, 1999.
- [Ram92] J. Ramanujam. Non-unimodular transformations of nested loops. In *Supercomputing 92*, pages 214–223. IEEE Computer Society Press, November 1992.
- [RS91] J. Ramanujam and P. Sadayappan. Compile-time techniques for data distribution in distributed memory machines. *IEEE Transactions on parallel and distributed systems*, 2(4):472–482, October 1991.
- [WR97] D. Wilde and S. Rajopadhye. Memory reuse analysis in the polyhedral model. *Parallel Processing Letter*, 7(2):203–215, 1997.