

Efficient Mapping of Interdependent Scans

Michel Barreteau and Paul Feautrier

Laboratoire PRiSM, Université de Versailles - S^t Quentin
45, avenue des États-Unis, 78 035 Versailles, FRANCE

Abstract. Distributed memory multiprocessors are extremely sensitive to communication costs. Some global communications such as scans and reductions are of special interest since their cost is much lower than for point to point communications. Our paper focuses on an algorithm which efficiently takes the mapping of scans into account.

1 Introduction

Communications remain the most critical aspect of performance in efficiently programming distributed memory multiprocessors. Hence minimizing communications is an indispensable task. A static placement may be obtained in two different ways: one may ask the user to insert annotations to specify data mapping à la HPF. Our approach is to leave this work to the compiler (see 2). However such a placement will not be sufficient if some particularities of the target machine are not taken into account, for instance if communication primitives with low overhead are not used. In this paper we propose a method to compute a placement which efficiently exploits interprocessor data movements such as scans and reductions. Especially we detail the Cholesky example to make the reader sensitive to the interdependences of scans. Neglecting them yields a mapping which may be incompatible with the minimization of communications.

2 Automatic parallelization

Automatic parallelization consists of extracting from the source program all restrictions on its potential parallelism. The Data Flow Graph (DFG) depicts the results of this analysis, i.e. the data movements between sets of operations [2]. A System of Affine Recurrence Equations (SARE) is built from the DFG. It is suitable for detecting recurrences [4]. We will consider the following SARE of three equations, which is extracted from a Cholesky decomposition program:

$$\begin{aligned}
 Ins_2[i, k] &= \begin{cases} Scan(\{i_2, k_2 \mid (3 \leq i_2 \leq n) \wedge (1 \leq k_2) \wedge (1 \leq i_2 - k_2)\}, \\ \quad ([0 \ 1]), +, -Ins_6[k_2, i_2]^2, InitS_2)[i, k] \\ \quad if \ (i \leq n) \wedge (2 \leq k) \wedge (1 \leq i - k) \\ a[i, i] - (a[k, i] / \sqrt{a[k, k]})^2 \quad if \ (2 \leq i \leq n) \wedge (k = 1) \end{cases} \\
 Ins_5[i, j, k] &= \begin{cases} Scan(\{i_5, j_5, k_5 \mid (3 \leq i_5) \wedge (j_5 \leq n) \wedge (1 \leq k_5) \\ \quad \wedge (1 \leq i_5 - k_5) \wedge (1 \leq j_5 - i_5)\}, ([0 \ 0 \ 1]), +, \\ \quad -(Ins_6[k_5, j_5] * Ins_6[k_5, i_5]), InitS_5)[i, j, k] \\ \quad if \ (j \leq n) \wedge (2 \leq k) \wedge (1 \leq i - k) \wedge (1 \leq j - i) \\ a[i, j] - ((a[k, j] / \sqrt{a[k, k]}) * (a[k, i] / \sqrt{a[k, k]})) \\ \quad if \ (2 \leq i) \wedge (j \leq n) \wedge (k = 1) \wedge (1 \leq j - i) \end{cases}
 \end{aligned}$$

$$Ins_6[i, j] = \begin{cases} Ins_5[i, j, i-1] / Ins_2[i, i-1] \\ \text{if } (2 \leq i) \wedge (j \leq n) \wedge (1 \leq j-i) \end{cases}$$

3 Scans & Reductions

Since scans and reductions are particular recurrences due to their associative operator, they may be computed in a parallel way using a binary tree. A reduction is a special scan in which only the final result is used. We only consider unidirectional scans, i.e. scans whose accumulation follows a unique direction in their iteration space. Let us remember the semantics of the unidirectional operator $\text{Scan}(\mathcal{D}_S, \mathbf{v}_S, \odot, \mathcal{D}_S, \mathbf{G}_S)$: \mathcal{D}_S describes the accumulation space, \mathbf{v}_S is the direction vector, \odot a binary and associative operator, \mathcal{D}_S the expression of the data to be reduced and \mathbf{G}_S are initial values. Scan computes a value for each integral point of \mathcal{D}_S , i.e. a multidimensional array \mathbf{S} . Some so-called *virtual* processors are in charge of these computations. They define the virtual processors space.

Two scans S_2 and S_5 whose direction vectors are $\mathbf{v}_{S_2} = (0, 1)$ and $\mathbf{v}_{S_5} = (0, 0, 1)$ have been detected in our SARE. S_2 computes the array $\mathbf{S2}$: $\forall (3 \leq i_2 \leq n) \wedge (1 \leq k_2) \wedge (1 \leq i_2 - k_2)$, $\mathbf{S2}[i_2, k_2] = \mathbf{S2}[i_2, k_2 - 1] + (-Ins_6[k_2, i_2]^2)$.

Practical implementations of the Scan operator allow only direction vectors which are parallel to the canonical axes of the virtual processors grid. A correct determination of the orientation of scans is necessary to avoid some additional communications related to non local data. But most of the time it depends on interactions between scans. Hence, the problem is to find out an efficient orientation of all the scans when placement binds them together. To this end we recall the successive stages of the mapping algorithm introduced in [3], which will be modified in order to solve our problem.

4 Placement algorithm

A placement function π_C , which is applied to a set of operations defined by an equation C of iteration vector \mathbf{i}_C , gives the name of the virtual processor on which C has to be executed. As regards geometry, a grid of virtual processors whose dimension g is given ($g = 2$ for Cholesky), has been elected. For each equation and for each scan, an affine placement function with g components has to be computed: for all dimension d ($1 \leq d \leq g$), $\pi_C^d(\mathbf{i}_C) = H_C^d \cdot \mathbf{i}_C + K_C^d \cdot \mathbf{n} + \mathbf{l}_C^d$ where unknowns are the coefficients μ_C of the $g \times |\mathbf{i}_C|$ matrix H , ν_C of the $g \times |\mathbf{n}|$ matrix K and ξ_C of the vector \mathbf{l} . The vector of structure parameters is \mathbf{n} (n belongs to it). $|\mathbf{v}|$ gives the number of components of a vector \mathbf{v} .

Ex: $\pi_{Ins_2}^d(i, k) = \mu_{Ins_2}^{d,1} \cdot i + \mu_{Ins_2}^{d,2} \cdot k + \nu_{Ins_2}^{d,1} \cdot n + \xi_{Ins_2}^d$

4.1 Placement conditions

First let us consider usual placement conditions that are given by the DFG:

$$\pi_{Ins_6}^d(i, j) = \pi_{Ins_5}^d(i, j, i-1) \Leftrightarrow \{\mu_{Ins_6}^{d,1} = \mu_{Ins_5}^{d,1} + \mu_{Ins_5}^{d,3}, \mu_{Ins_6}^{d,2} = \mu_{Ins_5}^{d,2}\} \quad (1)$$

$$\pi_{Ins_6}^d(i, j) = \pi_{Ins_2}^d(i, i-1) \Leftrightarrow \{\mu_{Ins_6}^{d,1} = \mu_{Ins_2}^{d,1} + \mu_{Ins_2}^{d,2}, \mu_{Ins_6}^{d,2} = 0\} \quad (2)$$

Then let us add the scan placement conditions which are developed in [1]:

1. Data alignment along the scan: $\forall \mathcal{C}(\mathbf{i}_C) \in \mathcal{D}_S(\mathbf{i}_S)$, $\pi_S(\mathbf{i}_S) = \pi_C(\mathbf{i}_C)$

$$\pi_{S_2}^d(i_2, k_2) = \pi_{Ins.6}^d(k_2, i_2) \Leftrightarrow \{\mu_{S_2}^{d,1} = \mu_{Ins.6}^{d,2}, \mu_{S_2}^{d,2} = \mu_{Ins.6}^{d,1}\} \quad (3)$$

$$\pi_{S_5}^d(i_5, j_5, k_5) = \pi_{Ins.6}^d(k_5, j_5) \Leftrightarrow \{\mu_{S_5}^{d,1} = 0, \mu_{S_5}^{d,2} = \mu_{Ins.6}^{d,2}, \mu_{S_5}^{d,3} = \mu_{Ins.6}^{d,1}\} \quad (4)$$

$$\pi_{S_5}^d(i_5, j_5, k_5) = \pi_{Ins.6}^d(k_5, i_5) \Leftrightarrow \{\mu_{S_5}^{d,1} = \mu_{Ins.6}^{d,2}, \mu_{S_5}^{d,2} = 0, \mu_{S_5}^{d,3} = \mu_{Ins.6}^{d,1}\} \quad (5)$$
2. Results collection from the scan to any variable \mathbf{X} such as $\mathbf{X}[\mathbf{i}_X] = \mathcal{S}[l(\mathbf{i}_X)]$ with $l(\mathbf{i}_X) = P_S \cdot \mathbf{i}_X + Q_S$: $\pi_X(\mathbf{i}_X) = \pi_S(l(\mathbf{i}_X))$

$$\pi_{Ins.2}^d(i, k) = \pi_{S_2}^d(i, k) \Leftrightarrow \{\mu_{Ins.2}^{d,1} = \mu_{S_2}^{d,1}, \mu_{Ins.2}^{d,2} = \mu_{S_2}^{d,2}\} \quad (6)$$

$$\pi_{Ins.5}^d(i, j, k) = \pi_{S_5}^d(i, j, k) \Leftrightarrow \{\mu_{Ins.5}^{d,1} = \mu_{S_5}^{d,1}, \mu_{Ins.5}^{d,2} = \mu_{S_5}^{d,2}, \mu_{Ins.5}^{d,3} = \mu_{S_5}^{d,3}\} \quad (7)$$

3. As regards the placement condition on computations which consists of finding for each scan an orientation following a j direction in the virtual processors grid: $\exists j$ ($1 \leq j \leq g$) such that $H_S^j \cdot \mathbf{v}_S \neq 0$ and $\forall i \neq j$, $H_S^i \cdot \mathbf{v}_S = 0$ (8), it should have priority to ensure an efficient computation. For each scan, g conditions $H_S^d \cdot \mathbf{v}_S = e_S^d$ are at the top of the system \mathcal{P} of placement equations in order to satisfy the strategy of the solution algorithm described below. The projection of \mathbf{v}_S by H_S in the virtual processors grid is \mathbf{e}_S . The components of \mathbf{e}_S are orientation parameters. All of them are equal to zero except one:

$\forall S$, $\mathbf{e}_S = k \cdot \mathbf{u}_S$ where $k \in \mathbb{Z}^*$ and \mathbf{u}_S is a canonical vector (9)

$$H_{S_2} \cdot \mathbf{v}_{S_2} = \mathbf{e}_{S_2} \Leftrightarrow \{\mu_{S_2}^{1,2} = e_{S_2}^1, \mu_{S_2}^{2,2} = e_{S_2}^2\} \quad (10)$$

$$H_{S_5} \cdot \mathbf{v}_{S_5} = \mathbf{e}_{S_5} \Leftrightarrow \{\mu_{S_5}^{1,3} = e_{S_5}^1, \mu_{S_5}^{2,3} = e_{S_5}^2\} \quad (11)$$

Most of the placement conditions have to be satisfied in order to minimize communications. A greedy algorithm based on the Gauss-Jordan elimination takes the edges that transfer the most important data volume first into account. This is their *weight*. In this way placement equations which are associated to these edges have a high probability to be cut and taken part in the solution of \mathcal{P} . (The order of edges to be considered according to their weight - as computed by our prototype - is: (7), (4), (5), (1), (3), (6) and (2).)

In order to avoid the mapping of every operation on a unique processor (experience shows that it often happens when all the equations are satisfied), a heuristic is adopted. This is the *triviality test*: it only accepts an equation if the current solution of \mathcal{P} (this equation included) is still able to generate g linearly independent solutions for each placement function.

The triviality test accepts for instance $\mu_{S_5}^{d,1} = 0$ from (4). But $\mu_{S_5}^{d,2} = 0$ is rejected because one cannot build for $Ins.5$ g (i.e. 2) linearly independent solutions any more ($\mu_{Ins.5}^{d,2} = \mu_{S_5}^{d,2}$ from (7) and $\mu_{Ins.5}^{d,1} = 0$ from (7) and (4)).

4.2 Interactions between scans

The pivot element of each equation is chosen to find out interactions between scans. As no assumption about relations on the coefficients μ_S can be made, we reject the problem to the relations \mathcal{S} between the orientation parameters e_S .

$\mu_{Ins.6}^{d,2} = \mu_{S_5}^{d,2}$, $\mu_{Ins.6}^{d,1} = \mu_{S_5}^{d,3}$ according to (4) and (3) implies $\mu_{Ins.6}^{d,2} = \mu_{S_2}^{d,1}$, $\mu_{Ins.6}^{d,1} = \mu_{S_2}^{d,2}$. Thus we get $\mu_{S_2}^{d,1} = \mu_{S_5}^{d,2}$ and $\mu_{S_2}^{d,2} = \mu_{S_5}^{d,3}$. Finally $e_{S_2}^d = e_{S_5}^d$ from (10) and (11).

The constraint to be satisfied for building \mathcal{S} comes straightforward from (8): let $(\mathbf{u}^1, \dots, \mathbf{u}^g)$ be the canonical vectors which describe the g directions of the virtual processors space, then \mathcal{S} must guarantee that each \mathbf{e}_S follows any canonical axis of the processors grid, i.e. each \mathbf{e}_S must be colinear to one \mathbf{u}^j .

Let us consider the equation \mathcal{E} : at first $e_{S_2}^1 = e_{S_5}^1$ since equations are processed dimension per dimension. Then a combination of canonical vectors \mathbf{u}_S^j which represents the direction of reorientation vectors \mathbf{e}_S is chosen. For instance $j = 1$ (cf (8)) for S_2 and S_5 , which means that both scans are supposed to lie in the direction of \mathbf{u}^1 . Rather than expressing this colinearity, it is easier to express that each \mathbf{e}_S has to be orthogonal to the $g - 1$ canonical vectors \mathbf{u}^j which describe the supplementary subspaces of the subspace defined by \mathbf{u}_S^j . In this way, one has to compute $g - 1$ dot products per scan: $\mathbf{e}_{S_2} \cdot \mathbf{u}^1 = 0$ and $\mathbf{e}_{S_5} \cdot \mathbf{u}^1 = 0$. Thus we get the following conditions $\{e_{S_2}^2 = 0, e_{S_5}^2 = 0\}$ which complete the subsystem $\mathcal{S}' = \mathcal{S} \cup \mathcal{E}$. (9) must be confirmed by the solution of \mathcal{S}' . $e_{S_5}^1 \neq 0$ is acceptable. This combination suits us. Otherwise we should choose another combination such as $(\mathbf{u}_{S_2}^1, \mathbf{u}_{S_5}^2)$. The second equation is: $e_{S_2}^2 = e_{S_5}^2$. The first combination $(\mathbf{u}_{S_2}^1, \mathbf{u}_{S_5}^1)$ yields: $e_{S_5}^2 = 0$. (9) is still verified. We deduce $\mathcal{S} = \mathcal{S} \cup \mathcal{E} = \{e_{S_2}^d = e_{S_5}^d\}$. On the other hand the second mentioned combination $(\mathbf{u}_{S_2}^1, \mathbf{u}_{S_5}^2)$ would yield: $\{e_{S_5}^1 = 0, e_{S_5}^2 = 0\}$. It would have been rejected because (9) is not satisfied any more. In this case \mathcal{S} would remain unchanged.

When every equation has been processed, one has to replace the parameters e_S^d by their own values into the global system \mathcal{P} using the last combination which satisfies \mathcal{S} . \mathcal{P} is solved and the d^{th} solution is applied to π_C^d . Thus we get:

$$\begin{array}{lll} \pi_{S_2}^1(i_2, k_2) = k_2 & \pi_{S_5}^1(i_5, j_5, k_5) = k_5 & \pi_{In_{S_5}}^1(i, j) = i \\ \pi_{S_2}^2(i_2, k_2) = i_2 & \pi_{S_5}^2(i_5, j_5, k_5) = j_5 & \pi_{In_{S_5}}^2(i, j) = j \end{array}$$

5 Conclusion

We presented in this paper an algorithm which computes multidimensional placement functions. It relies on data locality but also consider some special data movements such as scans and reductions and their interdependences. It enables to optimize mapping but also to produce an efficient target code. In this way hardwired communication primitives that are provided by many distributed memory multiprocessors will be efficiently exploited.

References

1. Barreteau, Feautrier: *Automatic Mapping of Scans and Reductions*. HPCS'95. 1995.
2. Feautrier: *Dataflow Analysis of Scalar and Array References*. Int. J. of Parallel Programming, 20(1):23-53. 1991.
3. Feautrier: *Toward Automatic Partitioning of Arrays on Distributed Memory Computers*. In ACM ICS'93 Tokyo pp.175-184. 1993.
4. Redon, Feautrier: *Detection of Reductions in Sequential Programs with Loops*. PARLE, LNCS 694. Ed. Arndt Bode and Mike Reeve and Gottfried Wolf. 1993.