



ELSEVIER

Parallel Computing 23 (1997) 251–266

PARALLEL
COMPUTING

Plugging anti and output dependence removal techniques into loop parallelization algorithm ¹

Pierre-Yves Calland ^{*},², Alain Darté, Yves Robert,
Frédéric Vivien

Laboratoire LIP, URA CNRS 1398, Ecole Normale Supérieure de Lyon, 69364 Lyon Cedex 07, France

Abstract

In this paper we shortly survey some loop transformation techniques which break anti or output dependences, or artificial cycles involving such ‘false’ dependences. These false dependences are removed through the introduction of temporary buffer arrays. Next we show how to plug these techniques into loop parallelization algorithms (such as Allen and Kennedy’s algorithm). The goal is to extract as many parallel loops as the intrinsic degree of parallelism of the nest authorizes, while avoiding a full memory expansion. We try to reduce the number of temporary arrays that we introduce, as well as their dimension.

Keywords: Anti dependence; Output dependence; Dependence removal; Loop parallelization algorithm

1. Introduction

Flow (or value-based) dependences are the only ‘true’ dependences of a program. Anti dependences ³ and output dependences ⁴ are due to storage re-use and can be eliminated at the price of more memory usage. Removing anti and output dependences

^{*} Corresponding author. E-mail: pierre-yves.calland@lip.ens-lyon.fr.

¹ Supported by the CNRS-ENS Lyon-INRIA Project *ReMaP* and by the Eureka project *EuroTOPS*.

² Supported by a grant of Région Rhône-Alpes.

³ Anti dependence occurs when a variable is used in one statement and reassigned in a subsequently executed statement [16].

⁴ Output dependence occurs when a variable is assigned in one statement and reassigned in a subsequently executed statement [16].

may prove very useful to break data dependence cycles and thereby enabling vectorization and/or improving parallelization.

However, removing *all* memory-based or ‘false’ (i.e., anti and output) dependences may have a prohibitive cost. A complete removal of false dependences is usually achieved, if feasible, via conversion of the original loop nest program into single assignment form. This turns out to be unnecessarily costly. Indeed, there are some memory-based dependences whose removal will not improve the parallelization. Rather, we should introduce as much memory overhead as needed to expose all the parallelism of the original program. As much as, but no more than, needed.

The aim of this paper is to show how to plug false dependence removal techniques into loop parallelization algorithms. The idea is to characterize those false dependences that do decrease the amount of parallelism, and to remove only these dependences. This will lead to memory savings without sacrificing performance.

The paper is organized as follows. Section 2 is devoted to a brief survey of techniques aimed at removing anti and output dependences. In Section 3 we work out an example to illustrate the key-ideas of our ‘integration’ sketch. We summarize the general steps of our method in Section 4. Finally, we give some conclusions in Section 5.

2. False dependence removal techniques

Many papers have been devoted to the problem of eliminating anti and output dependences. Proposed methods include ‘array data flow analysis’ [10,13], ‘variable expansion’ [4], ‘variable renaming’ [14] and ‘node splitting’ [14,6]. See the survey papers of Banerjee et al. [3] and Bacon et al. [2], as well as the books of Wolfe [16] and Zima and Chapman [17], for further references. Note that ‘array privatization’ [11] is yet another technique that can be applied, but it comes later, when moving from virtual to physical processors.

2.1. Renaming

Scalar renaming consists in giving a different name to occurrences of a scalar locally used in a program. This allows the removal of anti and output dependences due to the multiple use of the scalar. This technique can be directly extended to array renaming. Consider the loop nest in Fig. 1(a). The dependence graph⁵ (Fig. 1(c)) contains a cycle with an anti dependence from statement S_2 to statement S_3 . Renaming the array ‘a’ in S_3 (see the code in Fig. 1(b)) breaks this dependence: the new graph (Fig. 1(d)) is acyclic. Loop statements can now be parallelized.

⁵ In all figures, flow, anti and output dependence edges are labeled with a ‘f’, ‘a’ and ‘o’ respectively.

<p>For $i = 1$ to N do</p> <p>$S_1: a(i) = \sin(i)$</p> <p>$S_2: b(i + 1) = a(i) + c(i)$</p> <p>$S_3: a(i) = \cos(i)$</p> <p>$S_4: c(i + 1) = a(i)$</p> <p>EndFor</p> <p>(a) original code</p>	<p>For $i = 1$ to N do</p> <p>$S_1: \text{renamed}(i) = \sin(i)$</p> <p>$S_2: b(i + 1) = \text{renamed}(i) + c(i)$</p> <p>$S_3: a(i) = \cos(i)$</p> <p>$S_4: c(i + 1) = a(i)$</p> <p>EndFor</p> <p>(b) code after renaming</p>
--	--

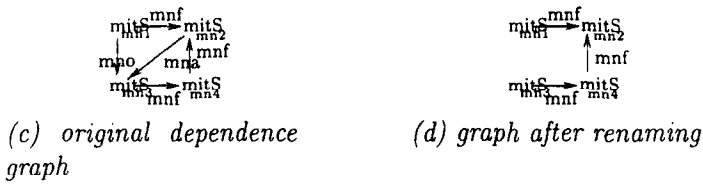


Fig. 1. Example of variable renaming.

<p>For $i = 1$ to N do</p> <p>$S_1: c(i) = 3 + a$</p> <p>$S_2: a = i + 1$</p> <p>$S_3: b(i) = c(i) + a$</p> <p>EndFor</p> <p>(a) original code</p>	<p>$\text{temp}(0) = a$</p> <p>For $i = 1$ to N do</p> <p>$S_1: c(i) = 3 + \text{temp}(i - 1)$</p> <p>$S_2: \text{temp}(i) = i + 1$</p> <p>$S_3: b(i) = c(i) + \text{temp}(i)$</p> <p>EndFor</p> <p>If $(N \geq 1)$ then $a = \text{temp}(N)$</p> <p>(b) code after expansion</p>
---	---

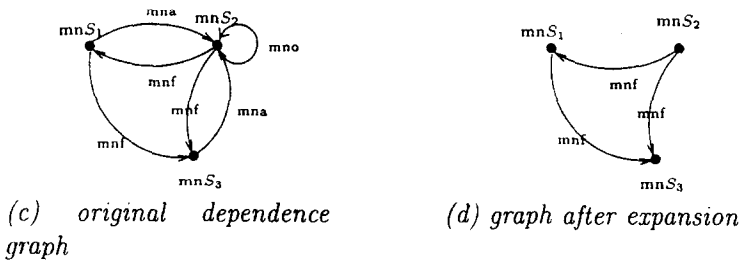


Fig. 2. Example of variable expansion.

2.2. Expansion

Consider a loop nest where a scalar variable is written at several iterations. This implies an output dependence from and to the statement involved in the multiple writings. Consider the example of Fig. 2(a). Since scalar a is written at each iteration, there is a self output loop around statement S_2 (see the dependence graph in Fig. 2(c)). To suppress this dependence, we expand a into a linear array, as shown in Fig. 2(b). Again, the new graph (Fig. 2(d)) is acyclic.

This technique can be extended to multi-dimensional loop nests for expanding scalars, or for expanding multi-dimensional arrays in the simple cases where the arrays can be considered as scalars when some loop indices are fixed.

2.3. Node splitting

This technique consists in splitting a statement into two statements, in order to break cycles in the dependence graph. Consider the example of Fig. 3(a). The dependence

For $i = 1$ to N do

$$S_1: a(i) = b(i) + c(i)$$

$$S_2: a(i + 1) = a(i) + 2 \times d(i)$$

(a) original code

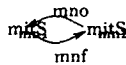
For $i = 1$ to N do

$$S'_1: temp(i) = b(i) + c(i)$$

$$S_1: a(i) = temp(i)$$

$$S_2: a(i + 1) = temp(i) + 2 \times d(i)$$

(b) code after node splitting



(c) original graph



(d) graph after node splitting

Fig. 3. Example of node splitting.

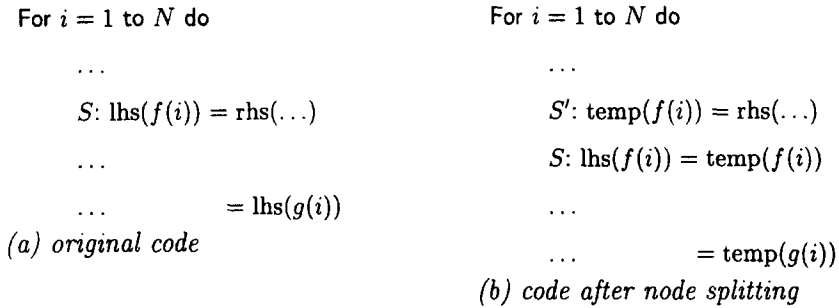


Fig. 4. Transformation of statement S .

graph (Fig. 3(c)) contains a cycle involving a flow dependence from S_1 to S_2 , and an output dependence from S_2 to S_1 . To break this cycle, rather than writing array ‘a’ in statement S_1 , we store the evaluation of the right-hand side into a temporary array ‘temp’. This temporary array is read in S_2 instead of array ‘a’. The transformed code is given in Fig. 3(b), and the new dependence graph is represented in Fig. 3(d).

The previous example is due to Padua and Wolfe [14]. We generalize [6] the statement transformation as indicated in Fig. 4. The value computed at each iteration of statement S is stored into a temporary array whose access function is the same as that of ‘lhs’, the left hand side of S . Obviously, if another statement instance depends upon a value ‘lhs($g(i)$)’ computed by S , then the access to ‘lhs($g(i)$)’ must be replaced by ‘temp($g(i)$)’. This implies knowledge of the statement instances which depend upon a value calculated in S (or in S' after the transformation).

The impact of this transformation on the dependences going to and coming from a statement S is summarized in Fig. 5. As shown in [6], if this transformation is applied to all the statements of the original loop nest, then the new dependence graph contains only flow dependence cycles and output dependence cycles. Furthermore, these cycles correspond to cycles of the initial dependence graph.

However, applying the transformation to all statements is not a good approach. First, it can be too costly. Moreover, it is useless to transform some statements. Consider for

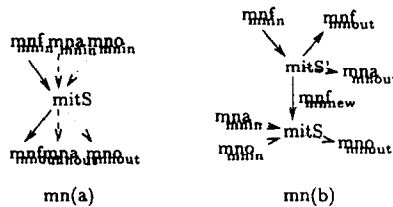


Fig. 5. A statement S with in-coming and out-going dependences (a) before and (b) after generalized Padua and Wolfe’s transformation.

instance a statement S (a vertex of the dependence graph) with an incoming flow dependence f_{in} and an outgoing output dependence o_{out} . We check on Fig. 5 that the path $\bullet \xrightarrow{f_{in}} S \xrightarrow{o_{out}} \bullet$ has been transformed into the path $\bullet \xrightarrow{f_{in}} S' \xrightarrow{f_{new}} S \xrightarrow{o_{out}} \bullet$. Thus a cycle containing the original path is not broken by the transformation. In fact, the paths that are broken when transforming statement S are those containing an anti or output dependence incoming to S , followed by an anti or flow dependence going out of S . To summarize, we break paths like

$$? \xrightarrow{a,o} S \xrightarrow{f,a} ?$$

2.4. Conversion to single-assignment form

The only full transformation to single assignment form (SAF) has been proposed by Feautrier in [10]. The technique relies on an exact analysis of direct flow dependences (through parametric integer linear programming) that permits the source of each array reference to be found, i.e. the statement and the value of the surrounding loop counters where the desired element of the array has been computed. Then, the algorithm is the following:

- For each statement S surrounded by d_S loops, define a new d_S -dimensional array M_S and replace the left-hand side of S by $M_S(I)$ where I is the iteration vector associated with S .
- Replace all references in the right-hand side using the corresponding source function: if the source is defined at iteration $f(I)$ of statement S' , then reference $M_S(f(I))$, and if the source is empty, then keep the original reference.

In other words, one temporary cell is introduced for each computation, and the right-hand side is modified to reference either a temporary cell for a value modified in the code, or the original scalar or array for a value kept unchanged. All new arrays have as many dimensions as there are loops surrounding their definition.

This transformation has two main weaknesses: the resulting code is in general very complicated, including many 'if' tests in the innermost loops; and it requires a very large amount of memory (one cell per computation). Some attempts have been made to remedy these two problems: more sophisticated rewriting techniques have been proposed to move 'if' tests into the outermost loops if possible and to minimize the memory usage (through memory folding, i.e. memory reuse) once parallelism has been detected (see the work of Chamski [7]).

In other words, Chamski's technique consists in three main steps:

1. transform the code into SAF, through full memory expansion,
2. parallelize the code,
3. reduce memory size by analyzing the life duration of each cell in the parallelized code.

In this paper, we explore an opposite approach: first determine anti and output dependences that are responsible for a loss of parallelism, remove them through memory expansion and then parallelize. We believe that this approach is more flexible and

powerful to enable various parallelization strategies. We point out that a similar approach is being currently developed in [12], but with a more restricted methodology since the false dependence removal is done *with respect to* a given schedule.

3. Motivating example

We briefly review Allen and Kennedy's algorithm (AK) in Section 3.1. Next we present a simple example, upon which we apply three parallelization schemes. First, in Section 3.2, we apply AK directly on the example. Then, in Section 3.3, we apply AK to the single assignment form of the example. Finally, in Section 3.4, we integrate the false dependences removal techniques into the parallelization process.

3.1. Allen and Kennedy's parallelization algorithm

We summarize Allen and Kennedy's algorithm (AK) because we use this parallelization algorithm throughout the paper. More details on this algorithm can be found in [1,5,17].

AK works on a structure called reduced leveled dependence graph (RLDG), i.e. a description of the level of dependences. For a *loop carried dependence*, the *level of dependence* is the rank of the first non null component of the distance vectors. This is also the depth of the outermost loop which carries this dependence. For a *loop independent dependence*, the level of dependence is said to be infinite and is denoted ∞ . If e is an edge of the RLDG, $l(e)$ denotes its level of dependence.

Before summarizing the algorithm in its simpler form, we need to recall some simple graph definitions:

- A *strongly connected component* of a directed graph G is a maximal subgraph of G in which for any vertices p and q ($p \neq q$) there is a path from p to q ;
- The *acyclic condensation* of a graph G is the acyclic graph whose nodes are the strongly connected components V_1, \dots, V_c of G . There is an edge from V_i to V_j if there is an edge $e = (x_i, y_j)$ in G such that $x_i \in V_i$ and $y_j \in V_j$.
- Let G be a reduced leveled dependence graph. Let H be a subgraph of G . Then $l(H)$ (the level of H) is the minimal level of an edge of H : $l(H) = \min\{l(e) | e \in H\}$.

$AK(H, l)$

- $H' = H \setminus \{e | l(e) < l\}$
 - Build H'' , the acyclic condensation of H' , and number its vertices V_1, \dots, V_c in a topological sort order.
 - For $i = 1$ to c do
 1. If V_i is reduced to a single statement S , with no edge, then generate parallel 'For' loops ('ForPar') in all remaining dimensions (i.e. for levels l to n) and generate code for S .
 2. Otherwise, let $k = l(V_i)$. Generate parallel 'For' loops ('ForPar') for levels from l to $k - 1$, and a sequential 'For' loop ('ForSeq') for level k . Call $AK(V_i, k + 1)$.
- Finally, to apply AK to a reduced leveled dependence graph G , call $AK(G, 1)$.

3.2. Direct parallelization scheme

The direct parallelization scheme consists in the application of AK on the following loop nest (with two 3-dimensional arrays a and c and one scalar b):

```

For  $i = 1$  to  $N$  do
  For  $j = 1$  to  $N$  do
    For  $k = 1$  to  $N$  do
       $S_1: a(i, j, k) = a(i, j, k + 1) + c(i, j + 1, k) + c(i, j - 1, k) + b$ 
       $S_2: b = a(i - 1, j, k) + a(i, j - 1, k)$ 
       $S_3: c(i, j, k) = c(i + 1, j, k) + c(i, j, k + 1)$ 
    EndFor
  EndFor
EndFor

```

Motivating example

The RLDG of the motivating example is drawn ⁶ on Fig. 6. This RLDG contains a single strongly connected component which includes the three statements and dependences at level 1. Thus, the outermost loop (loop i) is marked 'sequential' by AK. We now remove all level 1 edges: there is still a unique strongly connected component including at least one dependence at level 2. Thus, the second loop (loop j) is marked 'sequential'. We now remove level 2 edges: there are two strongly connected components, and each component includes dependences at level 3. Thus, the third loop (loop k) is marked 'sequential' for both components, and thus for all statements. Thus, AK finds no parallelism in this example when taking into account anti and output dependences, hence the need of removing at least some of the memory based dependences, in order to expose parallelism.

3.3. Parallelization of the single assignment form

Another approach could be first to transform the loop nest into single assignment form (thereby removing *all* memory based dependences), and then to apply the parallelization algorithm.

⁶ Dependences are those found by Tiny [15]. Some do not actually exist. For instance the anti dependence from S_1 to S_2 (due to scalar b) only occurs at level ∞ ; the anti dependences at level 1, 2 or 3 are covered.

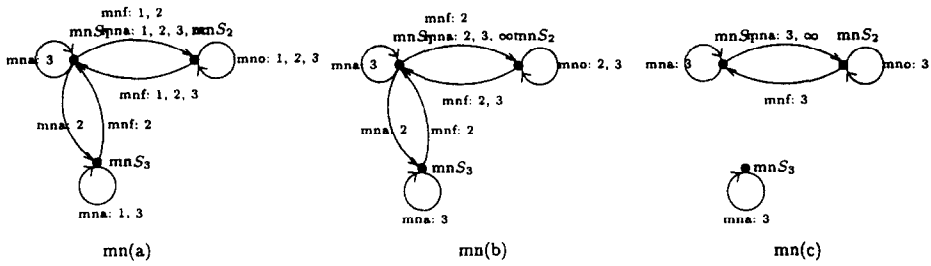


Fig. 6. RLDG of the motivating example: (a) whole RLDG; (b) RLDG without level 1 dependencies; (c) RLDG without level 1 and 2 dependencies.

Motivating example in single assignment form

We first transform the code into a single assignment form, using the transformations of Section 2.4 and we get (assuming $N \geq 1$):

```

For i = 1 to N do
  For j = 1 to N do
    For k = 1 to N do
      S1: atemp(i, j, k) = a(i, j, k + 1) + c(i, j + 1, k)
        + if j ≥ 2 then ctemp(i, j - 1, k) else c(i, 0, k)
        + if k ≥ 2 then btemp(i, j, k - 1)
          else if j ≥ 2 then btemp(i, j - 1, N)
            else if i ≥ 2 then btemp(i - 1, N, N)
              else b
      S2: btemp(i, j, k) = if i ≥ 2 then atemp(i - 1, j, k) else a(0, j, k)
        + if j ≥ 2 then atemp(i, j - 1, k) else a(i, 0, k)
      S3: ctemp(i, j, k) = c(i + 1, j, k) + c(i, j, k + 1)
    EndFor
  EndFor
EndFor
    
```

The new RLDG is drawn in Fig. 7.

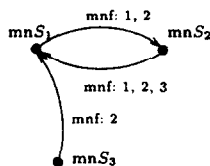


Fig. 7. RLDG of the motivating example in single assignment form.

Parallelization of the single assignment form

One can easily see that AK will mark the two outermost loops (loops i and j) 'sequential' for statements S_1 and S_2 . All other loops will be found 'parallel'. This is expressed in the parallelized form written below (strictly speaking, we should also copy back $\text{atemp}(i, j, k)$ into $a(i, j, k)$, $\text{ctemp}(i, j, k)$ into $c(i, j, k)$, and $\text{btemp}(N, N, N)$ into b).

```

ForPar  $i = 1$  to  $N$  do
  ForPar  $j = 1$  to  $N$  do
    ForPar  $k = 1$  to  $N$  do
       $S_3$ :  $\text{ctemp}(i, j, k) = c(i + 1, j, k) + c(i, j, k + 1)$ 
    EndFor
  EndFor
EndFor
ForSeq  $i = 1$  to  $N$  do
  ForSeq  $j = 1$  to  $N$  do
    ForPar  $k = 1$  to  $N$  do
       $S_2$ :  $\text{btemp}(i, j, k) = \text{if } i \geq 2 \text{ then } \text{atemp}(i - 1, j, k) \text{ else } a(0, j, k)$ 
        +  $\text{if } j \geq 2 \text{ then } \text{atemp}(i, j - 1, k) \text{ else } a(i, 0, k)$ 
    EndFor
    ForPar  $k = 1$  to  $N$  do
       $S_1$ :  $\text{atemp}(i, j, k) = a(i, j, k + 1) + c(i, j + 1, k)$ 
        +  $\text{if } j \geq 2 \text{ then } \text{ctemp}(i, j - 1, k) \text{ else } c(i, 0, k)$ 
        +  $\text{if } k \geq 2 \text{ then } \text{btemp}(i, j, k - 1)$ 
        else  $\text{if } j \geq 2 \text{ then } \text{btemp}(i, j - 1, N)$ 
          else  $\text{if } i \geq 2 \text{ then } \text{btemp}(i - 1, N, N)$ 
            else  $b$ 
    EndFor
  EndFor
EndFor

```

The latency of this parallel program is $\Theta(N^2)$, instead of $\Theta(N^3)$ for the direct parallelized version. Hence our motivating example does contain some parallelism! However, to expose all the parallelism, we have introduced three new arrays of size N^3 , which is the size of the iteration domain. We show in the next section that a clever integration of the false dependence removal techniques into the scheduling process enables maximum parallelism to be found while introducing less memory overhead.

3.4. Plugging false dependence removal techniques into the parallelization

The dataflow graph (the RLDG where only flow dependence edges are kept, see Fig. 8(a)) tells us exactly what amount of parallelism can be found in the program. Our aim is to find in the whole RLDG as much parallelism, while introducing as less memory overhead as possible.

3.4.1. First loop

Consider the first loop: because of the flow dependences, the outermost loop (loop i) must be sequential for statements S_1 and S_2 . However, this loop could be made parallel

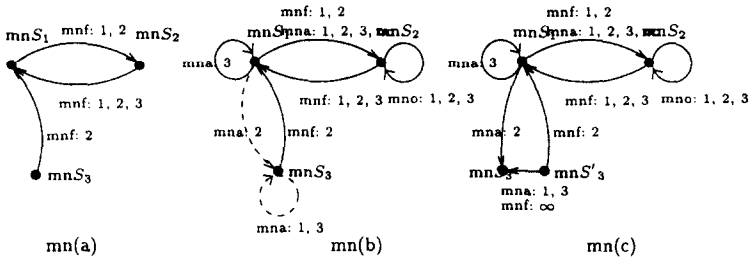


Fig. 8. Motivating example: (a) Dataflow graph (b) RLDG where incompatible edges are dashed (c) RLDG after program transformation.

for S_3 . In order to expose this parallelism, S_3 should not belong any longer to a strongly connected component including some level 1 dependence. Therefore two false dependence edges must be removed, we call them incompatible edges (see Fig. 8(b)):

- The anti dependence from S_1 to S_3 . Because of this edge and of the flow dependence from S_3 to S_1 , S_3 is in the same strongly connected component than S_1 and S_2 , which contains an edge at level 1.
- The self anti dependence on S_3 at level 1.

Note however that there is no need to remove the anti dependence of level 1 from S_1 to S_2 .

The two incompatible dependences can be removed by splitting the node S_3 as explained in Section 2.3. This only introduces a single new three dimensional array. The new RLDG is depicted in Fig. 8(c). We can now apply the first step of Allen and Kennedy's algorithm and we get:

```

ForPar i = 1 to N do
  ForPar j = 1 to N do
    ForPar k = 1 to N do
       $S_3$ :  $ctemp(i, j, k) = c(i + 1, j, k) + c(i, j, k + 1)$ 
    EndFor
  EndFor
EndFor
ForSeq i = 1 to N do
  For j = 1 to N do
    For k = 1 to N do
       $S_1$ :  $a(i, j, k) = a(i, j, k + 1) + b + c(i, j + 1, k)$ 
           + if  $j \geq 2$  then  $ctemp(i, j - 1, k)$  else  $c(i, 0, k)$ 
       $S_2$ :  $b = a(i - 1, j, k) + a(i, j - 1, k)$ 
    EndFor
  EndFor
EndFor
EndFor

```

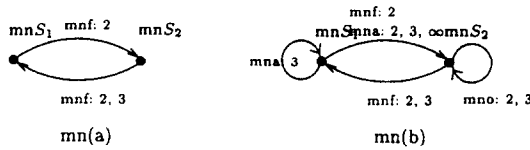


Fig. 9. After first level parallelization: (a) Dataflow graph (b) RLDG.

```

ForPar i = 1 to N do
  ForPar j = 1 to N do
    ForPar k = 1 to N do
      S3: c(i, j, k) = ctemp(i, j, k)
    EndFor
  EndFor
EndFor

```

3.4.2. Second loop

We now consider the second step of AK for the loops surrounding S_1 and S_2 (since the rest of the code is already fully parallelized). The remaining RLDG at level 2 is depicted in Fig. 9.

Because of the flow dependences at level 2 (see Fig. 9), the second loop (loop j) must be sequential for statements S_1 and S_2 . Removing the anti dependence at level 2 from S_1 to S_2 will not permit to detect more parallelism with AK. The second loop is marked sequential.

3.4.3. Third loop

Considering the dataflow graph of Fig. 10(a), we see that the innermost loop could be marked parallel because there is no cycle at level 3. In order to expose this parallelism, S_1 and S_2 should not belong any longer to a strongly connected component including some level 3 dependences. Therefore three false dependence edges must be removed (see Fig. 10(b)):

- the self output dependence on S_2 .
- the anti dependence from S_1 to S_2 .
- the self anti dependence on S_1 .

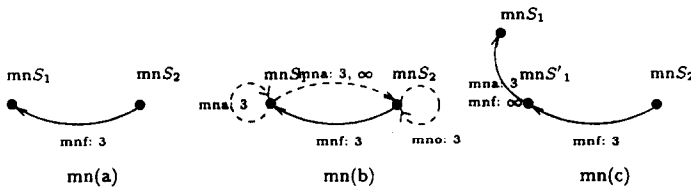


Fig. 10. After second level parallelization: (a) Dataflow graph (b) RLDG where incompatible edges are dashed (c) RLDG after program transformation.

The first two dependences can be removed by expanding the scalar b as explained in Section 2.2. The third dependence can be removed by splitting the node S_1 as explained in Section 2.3. However, instead of introducing two 3-dimensional arrays to suppress these dependences, we introduce only two 1-dimensional arrays ‘atemp’ and ‘btemp’: this is because the outermost two loops are already sequential. Indeed, we only need to remove incompatible dependences *for each iteration of the outermost loops*. This is done as if the outermost two loops indexes were fixed. We first get the code:

```

ForSeq i = 1 to N do
  ForSeq j = 1 to N do
    btemp(0) = b
    For k = 1 to N do
      S1: atemp(k) = a(i, j, k + 1) + btemp(k - 1) + c(i, j + 1, k)
        + if j ≥ 2 then ctemp(i, j - 1, k) else c(i, 0, k)
      S1: a(i, j, k) = atemp(k)
      S2: btemp(k) = a(i - 1, j, k) + a(i, j - 1, k)
    EndFor
    b = btemp(N)
  EndFor
EndFor

```

whose RLDG at level 3 is depicted in Fig. 10(c). Finally, applying the last step of AK leads to:

```

ForSeq i = 1 to N do
  ForSeq j = 1 to N do
    btemp(0) = b
    ForPar k = 1 to N do
      S2: btemp(k) = a(i - 1, j, k) + a(i, j - 1, k)
    EndFor
    ForPar k = 1 to N do
      S1: atemp(k) = a(i, j, k + 1) + btemp(k - 1) + c(i, j + 1, k)
        + if j ≥ 2 then ctemp(i, j - 1, k) else c(i, 0, k)
    EndFor
    ForPar k = 1 to N do
      S1: a(i, j, k) = atemp(k)
    EndFor
    b = btemp(N)
  EndFor
EndFor

```

4. Plugging false dependence removal techniques into parallelization algorithms

In Section 3, we have shown that integrating the false dependence removal techniques into the parallelization scheme makes it possible to find all the parallelism while introducing less memory overhead. In our example, one 3D array and two 1D arrays have been introduced instead of three 3D arrays as needed by the parallelization of the

single assignment form. Note that if external constraints had dictated that introducing a 3D array were too costly, we would have been able to cope with these constraints. We would have exposed less parallelism, but that is the price to pay!

We now summarize our methodology.

4.1. Integration scheme

Instead of removing all possible false dependences first and then applying a parallelization algorithm, we propose to combine both techniques.

Suppose that we use a parallelization algorithm that has the following properties:

- Each statement S surrounded by n_S loops in the original code is surrounded by n_S loops in the parallelized code.
- The choice of the loops surrounding S is made from the outermost to the innermost.

Note that all known parallelization algorithms have these properties: Allen and Kennedy's algorithm, Wolf and Lam's algorithm, Darte and Vivien's algorithm, Feautrier's algorithm.

We propose to plug false dependence removal techniques into loop parallelization algorithms as follows (where G is a given dependence graph):

Parallelization (G)

– Determine false dependences that could be removed by standard dependence removal techniques, as presented in Section 2. Let F be the dependence graph obtained if these dependences were removed. F is the dependence graph that exhibits as much parallelism as can be exposed.

– Apply the parallelization algorithm on F . Let PP be the parallelized program obtained. Each statement S is surrounded in PP by a sequence of n_S loops, marked either sequential or parallel.

– For $d = 1$ to $\max(n_S)$ do

(i) Mark as *incompatible* all dependences in G but not in F (i.e., false dependences that can be removed) which, if not removed at depth d , will induce a loss of parallelism if PP is taken as reference. In other words, incompatible dependences are those that prohibit the generation of the same number of parallel loops as in PP within the $n_S - d + 1$ remaining loops surrounding S .

(ii) Remove all the incompatible dependences by introducing temporary arrays of dimension as small as possible.

(iii) Generate the loop.

4.2. Comments

We do not go on more formally: the integration scheme above is simply the sketch of our methodology. Many problems remain to be solved: how to characterize incompatible edges for an arbitrary parallelization algorithm, how to minimize the number of incompatible dependences that should be removed, how to minimize the dimension of the temporary arrays that have to be introduced, ... However, we have given several examples in Section 3 which should make these problems clearer.

4.2.1. On the dimension of temporary arrays

In theory, we can hope to introduce temporary arrays with as many dimensions as nested loops minus the number of sequential loops already generated. See for example how node splitting has been performed in Fig. 10: we introduced only a 1D array and not a 3D array for splitting S_1 . This adds a new output dependence with level 1 and 2 for S_1' but it has no effect in terms of parallelization since the two outermost loops are already sequential.

In practice however, the dimension of the temporary arrays may be different. On one hand, we may need extra memory if the false dependence removal technique is not powerful enough to enable code generation: this is the case, for example, if the access functions are too complicated. On the other hand, memory overhead can be reduced by the use of scalar/array privatization, instead of expansion, along the parallel dimensions. This is illustrated in the example of Fig. 3: the 1D temporary array *temp* could be transformed into a privatized scalar.

4.2.2. Incompatible edges

Plugging dependence removal techniques into Allen and Kennedy's algorithm is straightforward, because it uses only loop distribution/fusion, which corresponds in terms of graph to the detection of strongly connected components in the RLDG. We have seen in Section 3 that this permits the identification of incompatible edges. They are precisely defined as follows:

For a RLDG H , we denote by H_l the subgraph of H obtained by deleting all edges with level $< l$. Then, incompatible edges at level l are the edges which belong to (at least) one cycle C such that:

1. C contains an edge of level l ,
2. C contains a vertex that only belongs in F_l to cycles of level strictly greater than l .

The characterization of incompatible edges for Darte and Vivien's algorithm is much more complicated. We refer to [8,9] for a complete description of the algorithm. We just give here the flavor of this algorithm.

Darte and Vivien's algorithm takes as input a reduced dependence graph G whose edges are labeled by dependence polyhedra. First, the reduced dependence graph G is uniformized into a graph G_u , which contains the nodes of G and some new nodes, called virtual nodes. Then G_u is processed by the parallelization algorithm as the reduced dependence graph of a system of uniform recurrence equations, except that one has to take into account the fact that some nodes are virtual. Then, the algorithm is recursive and it relies on the construction of a particular subgraph of G_u , the subgraph of null weight multi-cycles, denoted as G' . Incompatible edges are the edges that change the structure of G' , more precisely that change the structure of the vector space generated by the weight of the cycles of G' (or of one of the G' that will be recursively defined). We have not studied yet the complexity of this complete characterization.

5. Conclusion

We have presented a framework to plug false dependence removal techniques into loop parallelization algorithms. Our approach has two main advantages. First, we remove only false dependence edges that are responsible for a lesser degree of

parallelism, i.e. responsible for the sequentialization of an extra loop. Second, for each dependence removal — hence for each new temporary array — we use our knowledge on the already generated loops to (try to) minimize the number of dimensions of this temporary array.

Furthermore, our algorithmic sketch can be controlled by external parameters. For instance, we can require that each statement is surrounded by the same number of sequential loops. More important, we can straightforwardly cope with external constraints on the total memory overhead which is allowed. For instance if only 2D-arrays may be introduced, we will generate only 1D or 2D temporaries, at the price of some loss in parallelism.

Further work should be devoted to a better characterization of ‘incompatible’ edges, and to a precise estimation of the extra memory that is needed to expose all the potential parallelism.

References

- [1] J.R. Allen and K. Kennedy, Automatic translation of Fortran programs to vector form, *ACM Toplas* 9 (1987) 491–452.
- [2] D.F. Bacon, S.L. Graham and O.J. Sharp, Compiler transformations for high-performance computing, *ADM Comput. Surv.* 26(4) (1994).
- [3] U. Banerjee, R. Eigenmann, A. Nicolau and D.A. Padua, Automatic program parallelization, *Proc. IEEE* 81(2) (1993) 211–243.
- [4] T. Brandes, The importance of direct dependences for automatic parallelization, in: *Int. Conf. of Supercomputing* (1988) 407–417.
- [5] D. Callahan, A Global Approach to Detection of Parallelism, Ph.D. thesis, Dept. of Computer Science, Rice University, Houston, TX, 1987.
- [6] P.-Y. Calland, A. Darté, Y. Robert and F. Vivien, On the removal of anti and output dependences, in: J. Fortes, C. Mongenet, K. Parhi and V. Taylor, eds., *Application Specific Systems, Architectures and Processors* (IEEE Computer Society Press, 1996) 353–364.
- [7] Z. Chamski, Environnement logiciel de programmation d’un accélérateur de calcul parallèle, Ph.D. thesis, Université de Rennes, Rennes, France, 1993, No. 957.
- [8] A. Darté and F. Vivien, A classification of nested loops parallelization algorithms, in: *INRIA-IEEE Symp. on Emerging Technologies and Factory Automation* (IEEE Computer Society Press, 1995) 217–224.
- [9] A. Darté and F. Vivien, Optimal fine and medium grain parallelism detection in polyhedral reduced dependence graphs, Technical Report 96-06, LIP, ENS-Lyon, France, April 1996.
- [10] P. Feautrier, Dataflow analysis of array and scalar references, *Int. J. Parallel Program.* 20(1) (1991) 23–51.
- [11] J. Gu, Z. Li and G. Lee, Symbolic array dataflow analysis for array privatization and program parallelization, in: *Supercomputing 95* (1995).
- [12] V. Lefebvre and P. Feautrier, Gestion de la mémoire dans les programmes parallèles, in: R. Castanet and J. Roman, eds., *RenPur’8* (LaBRII, Université de Bordeaux, France, May 1996) 149–152 (in French).
- [13] D.E. Maydan, S.P. Amarasinghe and M. Lam, Array dataflow analysis and its use in array privatization, in: *Principles of Programming Languages* (1993).
- [14] D.A. Padua and M.J. Wolfe, Advanced compiler optimizations for supercomputers, *Commun. ACM* 29(12) (December 1986) 1184–1201.
- [15] M. Wolfe, The Tiny loop restructuring research tool, in: H.D. Schwetman, ed., *Int. Conf. on Parallel Processing*, Vol. II (CRC Press, 1991) 46–53.
- [16] M. Wolfe, *High Performance Compilers For Parallel Computing* (Addison-Wesley Publishing Company, 1996).
- [17] H. Zima and B. Chapman, *Supercompilers for Parallel and Vector Computers* (ACM Press, 1990).