

# On Index Set Splitting

Martin Griebel  
FMI, Universität Passau  
94030 Passau, Germany  
griebel@fmi.uni-passau.de

Paul Feautrier  
PRiSM, Université de Versailles  
45 avenue des États-Unis  
78035 Versailles, France  
Paul.Feautrier@prism.uvsq.fr

Christian Lengauer  
FMI, Universität Passau  
94030 Passau, Germany  
lengauer@fmi.uni-passau.de

## Abstract

*There are many algorithms for the space-time mapping of nested loops. Some of them even make the optimal choices within their framework. We propose a pre-processing phase for algorithms in the polytope model, which extends the model and yields space-time mappings whose schedule is, in some cases, orders of magnitude faster. These are cases in which the dependence graph has small irregularities. The basic idea is to split the iteration domain of the loop nests into parts with a regular dependence structure and apply the existing space-time mapping algorithms to these parts individually.*

*This work is based on a seminal idea in the more limited context of loop parallelization at the code level. We elevate the idea to the model level (our model is the polytope model), which increases its applicability by providing a clearer and wider range of choices at an acceptable analysis cost.*

*Index set splitting is one facet in the effort to extend the power of the polytope model and to enable the generation of competitive target code.*

## 1 Introduction

Space-time mapping methods for the automatic parallelization of loop nests relate every instance of every statement to a virtual point in time (*schedule*) and a virtual processor (*allocation*) [9, 14]. Linear algebra and linear programming are typically employed in the search for schedules and allocations. The information this search is based on is the set of dependences between different instances of the statements.

For the case of uniform dependences, Darte and Vivien [6] proved that there are methods which yield a schedule with (asymptotically) optimal latency. However, for the case of affine, non-uniform dependences, this

optimum is sometimes missed by orders of magnitude. The use of different schedules for different iterations of the same statement frequently improves this situation. Thus, the idea of this paper is, to partition the index sets of all statements independently of the problem size into a fixed number of parts and compute individual schedules for each part.

In the same way, it turns out that a piecewise defined placement function might improve the quality of an allocation, but this has never been pointed out in the literature.

## 2 Definitions

We are given a dependence graph whose nodes are statements and whose edges are dependences. Each node  $S$  is associated with a subset  $I(S)$  of  $\mathbb{N}^{p_S}$ , where  $p_S$  is the nesting depth of  $S$ . An edge  $e$  from  $S$  to  $T$  is associated with a dependence  $d_e$ , a relation from  $I(S)$  to  $I(T)$ .

The elements of  $I(S)$  are called *iteration vectors*. Each iteration vector is associated with an execution of statement  $S$ , which we also call an *operation*. The execution of  $S$  for iteration vector  $x$  is denoted  $\langle S, x \rangle$ . In more abstract contexts, operations are denoted by letters like  $u, v, \dots$ . The fact that  $v$  *depends* on  $u$  is written  $u \delta v$ . The set of operations, i.e., the disjoint union of all  $I(S)$  is named  $\Omega$ .

In the following, all iteration domains and dependence relations are parametric polytopes. In other words, they are defined by systems of affine constraints with parameters. To simplify the presentation, only one parameter, named  $n$ , is taken into account, meaning that the size of the iteration domain increases with  $n$ . We assume  $n$  to be unbounded.

A schedule is a function  $\theta$  from the set of operations to the set of integers which satisfies the following

causality condition:

$$\forall u, v \in \Omega : u \delta v \Rightarrow \theta(u) + 1 \leq \theta(v) \quad (1)$$

As a matter of fact, one can omit the integrity condition on schedules since, if  $\theta$  satisfies the causality condition, then so does  $\theta'(u) = \lfloor \theta(u) \rfloor$  [18].

From a causal schedule, one can deduce a parallel program whose figure of merit is its latency:

$$L = \max_{u \in \Omega} \theta(u) - \min_{u \in \Omega} \theta(u)$$

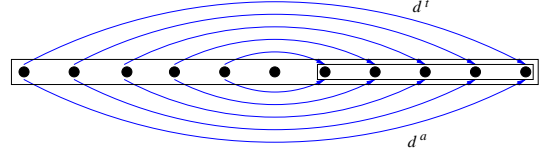
The latency can be interpreted either as the running time on a parallel computer with “enough” processors, or as the minimum number of synchronization points. Whatever the interpretation, it is clear that the latency must be minimized.

### 3 Statement of the Problem

It is not generally possible to find an arbitrary schedule with minimum latency. Usually, one restricts the search to a subset of all possible functions, the functions which are affine in the loop counters. Conceptually, one builds an affine template for every schedule function (i.e., an affine function with unknown coefficients) and writes inequality (1) for all possible values of  $u$  and  $v$ . The unknowns are the coefficients of the scheduling functions, and it is easy to see that the resulting constraints are affine in these unknowns. This must be done for all values of  $n$ , yielding an infinite set of constraints. Fortunately, thanks to special properties of affine functions, this set can be shown to be equivalent to a finite set of affine constraints, which can be solved by the usual linear programming methods.

For some problems, the resulting linear program is found to be infeasible. In this case, one resorts to multi-dimensional schedules. One can find a maximal subset of the dependences which still gives a feasible program. The resulting function is the first component of the multi-dimensional schedule. Then one applies the same algorithm to the unsatisfied dependences, obtaining the next component of the schedule, and so on until all dependences are satisfied.

One may wonder whether the schedule found in this way bears any relation to the minimum latency schedule. It has been proved [6] that, when all dependences are uniform, the two schedules are equivalent in the asymptotic sense. However, there are well known counter examples showing that this is not true for arbitrary affine dependences.



**Figure 1. Simple example showing the necessity of splitting**

**Example 1** Consider:

```
do i = 0, 2*n
  a(i) = a(2*n-i)
end do
```

The index set and its dependence graph are given in Figure 1. The best affine schedule is:

$$\theta_1(i) = i/2 \quad (2)$$

while the minimum latency schedule is:

$$\theta_2(i) = 0 \quad \text{if } 0 \leq i \leq n \quad (3)$$

$$= 1 \quad \text{if } n < i \leq 2n \quad (4)$$

$\theta_2$  can be found by splitting the iteration domain into two subdomains,  $I_1 = [0, n]$  and  $I_2 = [n+1, 2n]$ , and postulating two separate scheduling functions in  $I_1$  and  $I_2$ . The details of the resolution method are not affected by the splitting; the number of unknowns, however, is doubled. This splitting can also be interpreted as a code transformation yielding the program:

```
do i = 0, n
  a(i) = a(2*n-i)
end do
do i = n+1, 2*n
  a(i) = a(2*n-i)
end do
```

followed by the application of any convenient scheduling algorithm. ■

Our aim in this paper is to derive an algorithm for deciding when splitting is useful, and for finding these splits.

### 4 Related Work

Our notion of index set splitting seems very similar to tiling [2]: both techniques partition the index sets. Tiling is still a very active research area [10]. However, the goal of tiling has been either to increase granularity

(e.g., [3]), or to block for cache optimization (e.g., [5]), or simply to map virtual processors to real processors. In all these cases, the idea is to enumerate the given index set in a higher-dimensional space: one set of dimensions for the tiles and another set of dimensions for the points inside a tile; all tiles are treated equally. In contrast, index set splitting does not change the number of dimensions but benefits from an individual treatment of the various partitions.

Like index set splitting, the scheduling method by Feautrier [7, 8] can also result in piecewise affine functions. However, the schedules found are minima of a finite set of affine functions, and most piecewise affine schedules cannot be cast in this form. Example 1 is a case in point.

The idea of index set splitting goes back to Wolfe [19], and further to Allen/Kennedy [1] and Banerjee [4]. Our method expands on these seminal efforts by incorporating them into the polytope model.

The work most closely related is by Jemni and Mahjoub [15, 16] and deals with partitioning the index set at points where the type of a dependence changes, e.g., from true to anti. Since their method is not based on a model, they can separate the index domain only along planes parallel to the coordinate directions.

## 5 Analysis

Let us first explain why schedule (2) is suboptimal. An arbitrary affine function can be written as follows:

$$\theta(S, x) = \tau_S \cdot x + \sigma_S$$

and obviously has the property that

$$\theta(S, x + v) - \theta(S, x) = \tau_S \cdot v$$

depends on  $v$  but not on  $x$ . Suppose that there is a dependence from  $(S, x)$  to  $(S, x + v)$  for some  $x$  and  $v$ . Then:

$$\theta(S, x + v) - \theta(S, x) = \tau_S \cdot v \geq 1$$

Iterating this result  $k$  times, we obtain:

$$\theta(S, x + kv) - \theta(S, x) \geq k \quad (5)$$

The concept of latency can be extended to individual statements  $S$ :

$$L_S = \max_{x \in I(S)} \theta(S, x) - \min_{x \in I(S)} \theta(S, x)$$

and it is clear that the latency of a statement is a lower bound on the latency of the program. From (5) we deduce:

$$L_S \geq \max\{k \mid x + kv \in I(S)\} - \min\{k \mid x + kv \in I(S)\}$$

Since  $I(S)$  depends linearly on a size parameter  $n$ , we may expect that  $L_S$  also increases linearly with  $n$ .

In Example 1, we have a one-dimensional dependence vector,  $v = (2)$ , from instance  $n - 1$  to  $n + 1$ . Since  $v$  can be iterated  $n$  times within the index space, we have a latency of at least  $n$  for the execution of statement  $S$ , which is exactly the latency of schedule (2).

Suppose now that  $I(S)$  has been partitioned into two subdomains,  $I_1$  and  $I_2$ , and that  $x \in I_1$  and  $x + v \in I_2$ . The above reasoning no longer applies, since the schedule  $\theta$  is not necessarily the same affine function in  $I_1$  and  $I_2$ . Hence, the above estimate of the latency of  $S$  no longer holds, and we may hope for a better schedule.

If we split the index set of the target statement of a dependence  $d$  into  $I_1$ , which contains the image of  $d$  (and therefore has to satisfy the schedule constraints for  $d$ ), and  $I_2$ , which contains the rest of the index space (and therefore need not satisfy constraints which are due to  $d$ ), we get a constant schedule for  $I_2$  (if  $d$  is the only dependence). If there are no dependences inside  $I_1$ , i.e., the domain of  $d$  is contained in  $I_2$ , we also get a constant schedule for  $I_1$ .

For our example, the dependence analysis gives us the information that only the instances  $n+1, \dots, 2n$  are in the range of the existing dependences. Since iterations  $0, \dots, n$  are not images of any dependence, they need not satisfy any scheduling condition and, thus, should be given an individual schedule template. This splitting enables us to find schedule (3), which has a constant latency.

The same reasoning applies to placement problems. Here the goal is to find a placement function  $\pi(S, x) = \lambda_S \cdot x + \mu_S$  which gives the number of the processor that executes operation  $\langle S, x \rangle$ . The goal is that two operations which access the same memory location are executed on the same processor. If  $\langle S, x \rangle$  and  $\langle S, x + v \rangle$  access the same memory cell, we must have

$$\pi(S, c) = \pi(S, x + v) \Rightarrow \lambda_S \cdot v = 0$$

and all operations  $\langle X, x + kv \rangle$  will use the same processor as  $\langle S, x \rangle$ , whether there is a reason or not. On the other hand, putting  $\langle S, x \rangle$  and  $\langle S, x + v \rangle$  in different subsets of the iteration space invalidates this reasoning and may result in better parallelism. In the case of Example 1, the existence of a dependence vector  $v = (2)$  entails  $\lambda_S = 0$ , all operations are on the same processor, and there is no parallelism. After splitting, the only condition is that operations  $\langle S, x \rangle$  and  $\langle S, 2n - x \rangle$  are on the same processor, which can be arbitrary. The degree of parallelism is  $n$ .

## 6 A First, Naïve Splitting Algorithm

As we have seen in Section 5, suboptimal schedules result from the fact that some parts of a statement's index set are in the image of a dependence, and some are not. If the distinguished statement is the target of several dependences, we should apparently subdivide its iteration space, each part corresponding to a selection of the incoming dependences. Furthermore, if there is a dependence  $d_1$  from statement  $R$  to statement  $S$  and a dependence  $d_2$  from  $S$  to  $T$ , we should use the composite dependence  $d_2 \circ d_1$  for splitting  $I(T)$ . Thus, a naïve approach for index set splitting tries to compute all possible paths in the dependence graph and split every statement according to all incoming paths.

Even if some drawbacks are obvious, let us analyze this method in more detail. It will illustrate the limits of what we can expect from the final algorithm to be presented in Section 7.

### 6.1 Merging multiple incoming paths

First, we realize that a single statement can be reached via several paths. Thus, independently from the chosen splitting algorithm, we will have to merge the splits of multiple incoming dependences. In principle, we split the index set according to the first incoming dependence, then split every part according to the next incoming dependence, and so on.

It is important to see that the order of treating the incoming dependences is irrelevant: to merge splits means to compute a conjunction of the image (or its complement which we call the non-image) of an incoming dependence and the parts already split, and conjunction is associative and commutative.

The complexity of merging the splits of  $k$  arbitrary incoming (composite) dependences is  $2^k$ , i.e., in general, we must expect an algorithm with at least exponential time complexity. Therefore, our main concern must be to reduce the number  $k$  of different incoming (composite) dependences as much as possible, in order to get an efficient algorithm.

Note that, if the dependence graph is a tree, the complexity is only linear in the length of the paths, i.e., the depth of the tree: various incoming paths can only differ in length (for two different paths reaching a tree node, one must be a postfix of the other). Therefore, the image of the longer path is a subset of the image of the shorter path, i.e., we never need to split the non-images, which avoids the exponential growth.

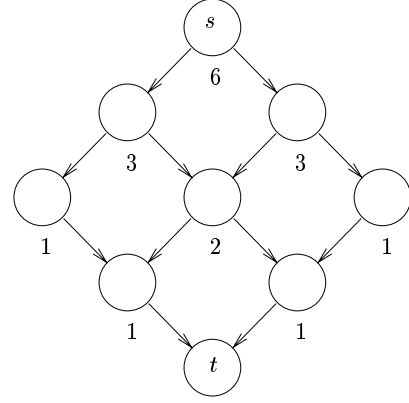


Figure 2. A dag with an exponential number of paths

### 6.2 Numbers of paths

We have just seen that in a tree the number of different paths to a given node is linear in the depth of the tree, i.e., logarithmic in the number of nodes. How many paths are there in a directed acyclic graph (DAG)?

There can be exponentially many paths between two statements in a DAG. For an illustration, consider Figure 2. For every node, we annotate the number of possible paths to node  $t$ . As we can see, the numbers correspond to Pascal's triangle, i.e., the numbers are the binomial coefficients  $\binom{l}{k}$ , where  $l$  is the length of the path from the considered node to  $t$ , and  $k$  is the horizontal position in the graph. In other words, the DAG in Figure 2 has  $O(l^2)$  nodes but  $O(2^l)$  paths.

Hence, if we consider all different paths from  $s$  to  $t$ , we end up with exponentially many incoming splits in  $t$ , which must then be merged with the exponential method as explained in Section 6.1. All in all, this results in a doubly exponential algorithm, which we consider impractical.

Obviously, the naïve method is not effective if there are loops in the statement dependence graph since, within strongly connected components, the number of different possible paths is unbounded. As simplest example, take a single statement with one self loop: every different number of loop traversals results in a different path. Note that, in this simple situation of only one self loop, the number of splits would increase linearly instead of exponentially (like in the case of the tree above), but it is still unbounded.

### 6.3 Preparing an effective algorithm

As we have just seen, directly using paths as descriptions of composite dependences results, in general, in a doubly exponential algorithm. Hence, our splitting algorithm must abstract from precise paths in order to be efficient.

For this purpose, we treat composite dependences systematically by associating a finite automaton with the dependence graph. The states of this automaton are the statements and the transitions are the dependences. There are well known algorithms (e.g., by Kleene [11]) for associating any two states  $S$  and  $T$  with a *regular expression* representing all paths from  $S$  to  $T$ . The letters in this regular expression are dependence *names*, and the operators are the dot (concatenation), the vertical bar (set union), and the Kleene star. The composite dependence from  $S$  to  $T$  is obtained by replacing in this expression each dependence name by the dependence itself, the dot by relation composition, the vertical bar by relation union, and the Kleene star by transitive closure. In suitable cases, one can compute the composite dependence in closed form by using, for instance, the Omega calculator [17]. However, since the transitive closure of an affine relation is not always affine, the above computation does not always succeed. Our proposal is to ignore a composite dependence when a closed form cannot be computed.

However, if we base our splitting algorithm on the path descriptions given by the Kleene algorithm, we find many practical examples which cannot but should be split. This is because the use of the operators  $|$  and  $*$  entails a loss of precision. In fact, we lose all information on the *delay* associated with the composite dependence.

**Example 2** Consider the following program:

```
do i = 0, m
  do j = 0, 2*n
    a(i,j) = a(i,2*n-j) + a(i-1,j+2)
  end do
end do
```

The iteration dependence graph is depicted in Figure 3. The range of the combined dependence  $d = d_{\text{upwards}} \mid d_{\text{downwards}}$  is  $[0, m] \times [0, 2n] \setminus [0, 0] \times [0, n]$ , as is the range of  $d^+$ . The desired split into  $[0, m] \times [0, n]$  and  $[0, m] \times [n+1, 2n]$ , which leads to a linear execution time (see Figure 3), is not found if we base the algorithm on Kleene's path descriptions only. ■

The difficulty in this example is that the different dependences have individual properties, which are merged

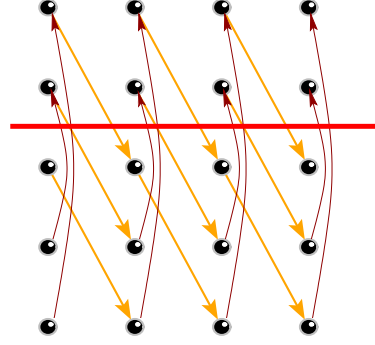


Figure 3. Splitting due to the initial phase

by the union operation, making the distinctions invisible. However, treating all dependences separately is too costly, in general, as shown in Section 6.2.

Our heuristic solution is to consider (in addition to all combined paths) every single dependence once, and split the index set of its target statement into its range and the rest. This guarantees that the properties of different dependences are considered at least once, and that complexity is increased by only a linear factor ( $D$  dependences instead of  $2^D$  paths). In other words, by splitting the target index set  $I(S)$  of every single dependence  $d$  into the range of  $d$  and the rest, we already have a conservative approximation of the range of all *paths* to  $S$  whose last dependence is  $d$ .

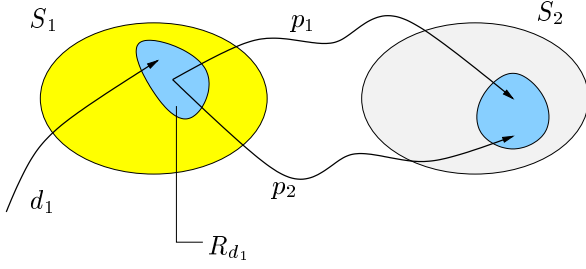
Our idea is then to propagate these ranges of  $d$  along every path to all other statements, where we now accept the loss of information by using the path descriptions from Kleene's algorithm, in order to keep the computational effort reasonably small.

Practical experiments show that this mixture of working with separate dependences on the one hand and combined path descriptions on the other exhibits a good balance between power and execution complexity of our algorithm.

## 7 The Proposed Splitting Algorithm

With the ideas in the previous section, we can now formulate a first version of the effective splitting algorithm:

1. For all dependences  $d$ , compute a polytope  $R_d$  (as small as possible) containing the range of  $d$ , and split the index set  $I(T)$  of the target statement of  $d$  into  $R_d$  and  $I(T) \setminus R_d$ .
2. Compute a description for the set of all paths in



**Figure 4. An illustration of the splitting algorithm**

the statement dependence graph, using Kleene's algorithm.

3. For every pair of statements  $(T, S)$  and for every dependence  $d$  with target statement  $T$  and every path  $p$  from  $T$  to  $S$  do: interpret path description  $p$  as a composition of relations which maps points of the index space  $I(T)$  to points of the index space  $I(S)$ , and compute the image  $p(R_d)$  of this composed relation when applied to the polytope  $R_d$  computed in Step 1. This will divide the index set  $I(S)$  into a part which is in the image of  $R_d$  under  $p$ , and the rest. Usually, this step can be computed with the Omega calculator [17]. However, if  $p$  contains a cycle whose transitive closure cannot be computed precisely by Omega, then delete the cycle from  $p$  before the propagation.
4. For any statement  $S$ , combine all splits obtained in this way as in Section 6.1.

The algorithm is illustrated in Figure 4. In Step 1, the index set of  $S_1$  is split into the image  $R_{d_1}$  of  $d_1$  and the rest. Step 2 computes all paths between every pair of statements. Assume that there are two paths from  $S_1$  to  $S_2$ , denoted with  $p_1$  and  $p_2$ . Step 3 computes the image of  $R_{d_1}$  w.r.t.  $p_1|p_2$  within the index space of  $S_2$ , where the possible images are all in the dark subset.

Basically, this algorithm computes, for every statement, the approximate ranges of all (transitively) incoming dependences. It is obvious that, for each such range, we might obtain a different possible schedule and, thus, want a separate template.

However, in practice, there are some additional considerations:

- It is easy to see that splitting an index set due to a uniform dependence is worthless, since the range of a uniform dependence is (approximately) the complete target index space. Therefore, we

optimize the algorithm by applying Step 1 only to non-uniform dependences. Note that, in Step 3, we still need to consider uniform dependences, as we shall see in Example 4.

- Due to the condensed description of the set of all paths and the overestimation of the reflexive transitive closure by Omega, it often happens that we lose precision and, thus, do not find precise splits and sometimes even not all splits (if two approximations yield the same set).

Note that these overestimates by Omega are due partly to the theoretical impossibility of computing the reflexive transitive closure as a finite union of polyhedra and partly to technical limitations inside Omega.

- Sometimes it is useful to unroll a cycle a fixed number of times, in order to achieve the optimal split (cf. Example 6). However, for complexity reasons, we decided not to take this idea into account in our prototype implementation of the method. An analysis of when unrolling is useful in general is left for future work.

## 8 Examples

Let us reconsider our initial example and modify it in several ways in order to get a feeling for the performance of our method.

**Example 3** Let us apply our algorithm to Example 1. The iteration domain is the set  $[0, 2n]$  and the range  $R$  of the two dependences is  $[n+1, 2n]$ . The set of all paths in the statement dependence graph is given by  $(d_t | d_a)^*$ . Propagating  $R$  along  $d^+$  (+ meaning at least once) is not possible since the domain of  $d^+$  does not intersect  $R$ . Thus, the algorithm already terminates after the initial step and splits the index set into  $[0, n]$  and  $[n+1, 2n]$ , as expected. ■

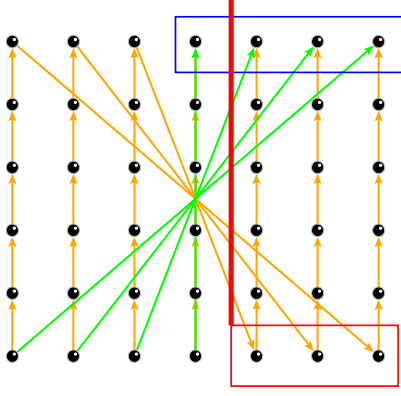
**Example 4** Extending the program of Example 1 to a two-dimensional example, let us consider:

```

do i = 0, 2*n
  do j = 0, m
    a(i, j) = a(2*n-i, j+m) + a(i, j-1)
  end do
end do

```

For the uniform dependence, no split is derived. The non-uniform true dependence has range  $[n+1, 2n] \times [0, 0]$ . Propagating along the combined self-dependence



**Figure 5. Splitting a two-dimensional index set by propagation**

leads to the desired split into  $[0, n] \times [0, m]$  and  $[n+1, 2n] \times [0, m]$  (see Figure 5). The non-uniform anti-dependence has the range  $[n, 2n] \times [m, m]$ , which is not increased by propagation. So, finally, we end up with three subdomains:  $[0, n] \times [0, m] \setminus \{(n, m)\}$ , the singleton  $\{(n, m)\}$ , and  $[n+1, 2n] \times [0, m]$ . ■

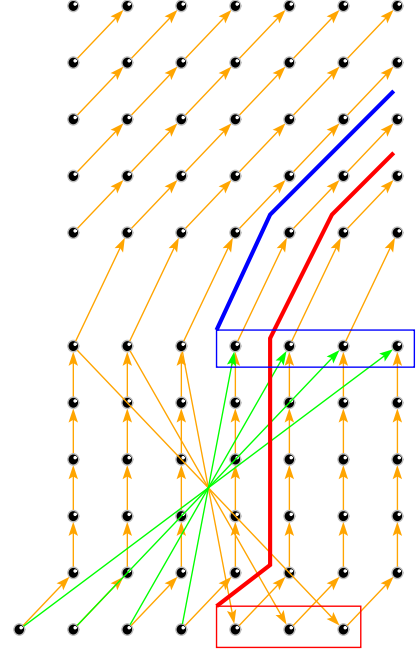
Note that, in the previous example, treating the singleton individually does not improve the schedule and, hence, should be avoided in practical implementations, if possible. However, in general, it is impossible to decide locally, i.e., at one given statement, whether a split is useful or not, since global information of the dependence graph is necessary. In other words, we would have to run the scheduler in order to detect whether we should split the input of the scheduler! Thus, in our current implementation, we accept possibly useless splits.

**Example 5** In order to get a better feeling for the behaviour of our algorithm in the case of multiple statements and imperfect loop nests, let us consider:

```

do i = 0, 2*n
S1:  a(i,0) = a(2*n-i,m)
    do j = 1, m
S2:    a(i,j) = a(i,j-1)
    end do
end do
do i = 0, 2*n
    do j = 1, k
S3:    a(i,j+m) = a(i-1,j+m-1)
    end do
end do

```



**Figure 6. Splitting for imperfect loop nests**

In this example, we find two uniform true dependences and four non-uniform dependences (three true and one anti), i.e., we obtain four initial splits. After propagating each of them along the combined path to each of the three statements, we have to merge the four initial splits with the 12 propagated splits. After simplification, we get the desired splits as indicated in Figure 6. ■

**Example 6** In the case of a self-dependence  $d$ , a constant number of propagations may give us an interesting schedule. Let  $I$  be the iteration space of the dependence. Successive propagation yields subsets  $I - d(I), d(I) - d^2(I), \dots, d^r(I) - d^{r+1}(I)$  and, in the interesting case, there exists a constant  $k$  such that  $d^k = \emptyset$ . It can be proved that  $k$  is bounded by the cardinality of  $I$ , but this bound depends on the size parameters of the program and is not constant. Consider the following program:

```

do i = 0, 5*n-1
  a(i) = a(i-n)
end do

```

If we treated the only existing dependence  $d$  as quasi-uniform, we would not split at all; treating  $d$  as non-uniform, we split the index set initially into  $[0, n-1]$  and  $[n, 5n-1]$ . Since the image of  $[n, 5n-1]$  under  $d^*$

is  $[n, 5n-1]$ , we end up with the two partitions already given. However, in this case, it is easy to prove that  $k = 5$ , and hence that the latency of the program is 5 instead of  $4n + 1$ .

We do not know at present whether there exists an algorithm for finding  $k$ , or whether the problem is undecidable. This question is left for future work.

Another approach to this example would be to try to find a scheduler which yields  $\theta(S) = \lfloor i/n \rfloor$ . Formally, this schedule has linear execution time but, in practice, it is constant (five steps) if we take the loop bounds into account. Unfortunately, such a solution cannot be found in the affine mathematical setting; only if  $n$  were a constant number instead of a parameter, existing schedulers could find this solution [7]. It is left for future work to check whether existing schedulers can be given the capability of symbolic computations, which would allow them to solve this example optimally. ■

## 9 Conclusions

We have proposed a method to split the iteration domain of loop nests into parts in order to obtain better space-time mappings. It can be used to improve any space-time mapping algorithm which is based on templates (templates are typically used, e.g., for dealing with different statements in the loop body).

Our method can be adapted simply for trading off quality of parallelism for analysis time:

- One can search for more splits and, thus, even satisfy a variety of optimality criteria, but at the cost of a doubly exponential search.
- One can turn off the propagation phase in order to save compilation time at the price of a loss of some useful splits (e.g., in Example 4).
- Since the proposed method is quite costly, it would be very interesting to be able to apply it where and only where it is potentially useful. A suggestion is to iterate between scheduling and splitting, but apply this technique only to statements whose schedule is not satisfactory. This topic is further considered in [12].

In our experiments, we decided to use the presented algorithm, which seems to be a good compromise between analysis time and quality of the resulting splits. More experiments on practical examples with our prototype will have to confirm this observation.

For this purpose, we are currently implementing the described methods in the prototype parallelizer LooPo [13]. This implementation takes as input the results

of the existing dependence analysis modules and is itself very close to the algorithm in Section 7, with the following technical modifications:

- In Step 1, we do not actually execute the splits but just store them per statement in a list, in order not to increase the number of statements in this phase.
- In Step 3, we only compute the images; the computation of the rest (i.e., the non-images) does not commence before the end of Step 4 and, again, the derived splits are just stored in lists in order to limit complexity.
- In Step 4, we first merge all stored splits and then compute the rest, i.e., the non-images – which may lead to further splitting, due to technical limitations, if the non-images cannot be described by a single polytope but only by a union of polytopes. Finally, every split, i.e., description of a subset of the original index space gets a copy of the body statements and, thus, can be viewed by all subsequent parallelization phases as if it were a separate statement with surrounding loops in the input program. Note that our model-based approach saves us from computing a loop nest which really enumerates all these split subsets; this is advantageous because the construction of such a loop nest would be very costly.

More details on unrolling self-dependences, as illustrated in Example 6, and targetting specific statements for a split can be found in an expanded version of this paper [12].

## Acknowledgements

Financial support was gratefully received from the German Research Foundation (DFG) under project *RecuR* and from the French-German exchange programme *PROCOPE* through DAAD and APAPE.

Thanks to Mohamed Jemni for discussing his method with us in Passau, which started off this work. We are grateful to Jean-François Collard and Bernhard Lehner for fruitful discussions.

## References

- [1] John R. Allen and Ken Kennedy. Automatic translation of FORTRAN programs to vector form. *ACM Trans. on Programming Languages and Systems*, 9(4):491–542, October 1987.



- [2] Corinne Ancourt and François Irigoin. Scanning polyhedra with DO loops. In *Proc. 3rd ACM SIGPLAN Symp. on Principles & Practice of Parallel Programming (PPoPP'91)*, pages 39–50. ACM Press, 1991.
- [3] Ruman Andonov, Sanjay Rajopadhye, and Nicola Yanev. Optimal Orthogonal Tiling. In David Pritchard and Jeff Reeve, editors, *Euro-Par'98: Parallel Processing*, LNCS 1470, pages 480–490. Springer-Verlag, 1998.
- [4] Utpal Banerjee. *Speedup of Ordinary Programs*. PhD thesis, Department of Computer Science, University of Illinois at Urbana-Champaign, October 1979. Report 79-989.
- [5] Rajeev Barua, David Kranz, and Anant Agarwal. Communication-minimal partitioning of parallel loops and data arrays for cache-coherent distributed-memory multiprocessors. In David Sehr, Utpal Banerjee, David Gelernter, Alex Nicolau, and David Padua, editors, *Languages and Compilers for Parallel Computing (LCPC'96)*, LNCS 1239, pages 350–368. Springer-Verlag, 1997.
- [6] Alain Darte and Frédéric Vivien. On the optimality of Allen and Kennedy's algorithm for parallelism extraction in nested loops. In Luc Bougé, Pierre Fraigniaud, Anne Mignotte, and Yves Robert, editors, *Euro-Par'96: Parallel Processing, Vol. I*, LNCS 1123, pages 379–388. Springer-Verlag, 1996.
- [7] Paul Feautrier. Some efficient solutions to the affine scheduling problem. Part I. One-dimensional time. *Int. J. Parallel Programming*, 21(5):313–348, 1992.
- [8] Paul Feautrier. Some efficient solutions to the affine scheduling problem. Part II. Multidimensional time. *Int. J. Parallel Programming*, 21(6):389–420, 1992.
- [9] Paul Feautrier. Automatic parallelization in the polytope model. In Guy-René Perrin and Alain Darte, editors, *The Data Parallel Programming Model*, LNCS 1132, pages 79–103. Springer-Verlag, 1996.
- [10] Jeanne Ferrante, Wolfgang Giloi, Sanjay Rajopadhye, and Lothar Thiele, editors. Tiling for optimal resource utilization. Technical Report 221, Schloß Dagstuhl, August 1998.
- [11] Robert W. Floyd and Richard Beigel. *The Language of Machines – An Introduction to Computability and Formal Languages*, chapter 4.4. Computer Science Press, 1994.
- [12] Martin Griebl, Paul F. Feautrier, and Christian Lengauer. Index set splitting. Technical Report MIP-9908, Fakultät für Mathematik und Informatik, Universität Passau, August 1999.
- [13] Martin Griebl and Christian Lengauer. The loop parallelizer LooPo—Announcement. In David Sehr, Utpal Banerjee, David Gelernter, Alex Nicolau, and David Padua, editors, *Languages and Compilers for Parallel Computing (LCPC'96)*, LNCS 1239, pages 603–604. Springer-Verlag, 1997.
- [14] Christian Lengauer. Loop parallelization in the polytope model. In Eike Best, editor, *CONCUR'93*, LNCS 715, pages 398–416. Springer-Verlag, 1993.
- [15] Zaher Mahjoub and Mohamed Jemni. Restructuring and parallelizing a static conditional loop. *Parallel Computing*, 21(2):339–347, February 1995.
- [16] Zaher Mahjoub and Mohamed Jemni. On the parallelization of single dynamic conditional loops. *Simulation Practice and Theory*, 4:141–154, 1996.
- [17] William Pugh and Dave Wonnacott. Eliminating false data dependences using the Omega test. *ACM SIGPLAN Notices*, 27(7):140–151, July 1992. *Proc. ACM SIGPLAN Conf. on Programming Language Design and Implementation (PLDI'92)*.
- [18] Patrice Quinton. The systematic design of systolic arrays. In Françoise F. Soulié, Yves Robert, and Maurice Tchuenté, editors, *Automata Networks in Computer Science*, chapter 9, pages 229–260. Manchester University Press, 1987. Also: Technical Reports 193 and 216, IRISA (INRIA-Rennes), 1983.
- [19] Michael Wolfe. *Optimizing Supercompilers for Supercomputers*. Research Monographs in Parallel and Distributed Computing. MIT Press, 1989.