



University of Delaware  
Department of Electrical and Computer Engineering  
Computer Architecture and Parallel Systems Laboratory

---

## Single-Dimension Software Pipelining for Multi-Dimensional Loops

*Hongbo Rong*†  
*Zhizhong Tang*††  
*R. Govindarajan*†  
*Alban Douillet*†††  
*Guang R. Gao*†

**CAPSL Technical Memo 049**

Revised on January 7, 2007

Copyright © 2005 CAPSL at the University of Delaware

†Dept. of Electrical and Computer Engineering  
University of Delaware  
rong,govind,ggao@capsl.udel.edu  
††Dept. of Computer Sciences  
University of Delaware  
douillet@capsl.udel.edu  
††Dept. of Computer Science and Technology  
Tsinghua University  
tzz-dcs@tsinghua.edu.cn



## Abstract

Traditionally, software pipelining is applied either to the innermost loop of a given loop nest or from the innermost loop to outer loops. This paper proposes a 3-step approach, called *Single-dimension Software Pipelining (SSP)*, to software pipeline a loop nest at an arbitrary loop level that has a rectangular iteration space and contains no sibling inner loops in it.

The first step identifies the most profitable loop level for software pipelining in terms of initiation rate, data reuse potential, or any other optimization criteria. The second step simplifies the multi-dimensional data-dependence graph (DDG) of the selected loop level into a 1-dimensional DDG and constructs a 1-dimensional schedule. Based on it, the third step derives a simple mapping function that specifies the schedule time for the operation instances in the multi-dimensional loop.

The classical modulo scheduling is subsumed by SSP as a special case. SSP is also closely related to hyperplane scheduling, and, in fact, extends it to be resource-constrained. We prove that SSP schedules are correct, and at least as efficient as those schedules generated by traditional modulo scheduling methods.

We extend SSP to schedule imperfect loop nests, which are most common at instruction-level. Multiple initiation intervals are naturally allowed to improve execution efficiency.

Feasibility and correctness of our approach are verified by a prototype implementation in the ORC compiler for the IA-64 architecture, tested with loop nests from Livermore and SPEC2000 floating-point benchmarks. Preliminary experimental results reveal that, compared to modulo scheduling, software pipelining at an appropriate loop level results in significant performance improvement. Software pipelining is beneficial even with loop transformations beforehand.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Basic Concepts</b>	<b>4</b>
<b>3</b>	<b>Motivation and Overview</b>	<b>5</b>
3.1	A Motivating Example . . . . .	5
3.2	Illustration of Our Approach . . . . .	6
3.3	Overview of Our Approach . . . . .	9
3.3.1	Which Loop to Software Pipeline? . . . . .	9
3.3.2	How to Software Pipeline the Selected Loop? . . . . .	9
3.3.3	Constructing the Final Schedule . . . . .	10
<b>4</b>	<b>Solution</b>	<b>10</b>
4.1	Definition of Simplified DDG . . . . .	11
4.2	Step 1: Loop Selection . . . . .	13
4.2.1	Initiation Rate . . . . .	13
4.2.2	Data Reuse . . . . .	14
4.3	Step 2: 1-D Schedule Construction . . . . .	15
4.4	Step 3: Final Schedule Computation . . . . .	16
<b>5</b>	<b>Analysis</b>	<b>17</b>
5.1	Correctness . . . . .	17
5.2	Efficiency . . . . .	18
5.3	Relation to the Classical Modulo Scheduling of a Single Loop . . . . .	20
5.4	Relation to Hyperplane Scheduling . . . . .	20
5.5	Time Complexity . . . . .	21
<b>6</b>	<b>Extension to Imperfect Loop Nests</b>	<b>21</b>
6.1	Motivating Example . . . . .	22
6.2	Assumptions . . . . .	23
6.3	Solution with a Single Initiation Interval . . . . .	23
6.3.1	1-D Schedule Construction . . . . .	25
6.3.2	Final Schedule Computation . . . . .	27
6.3.3	Relation to the Scheduling of Perfect Loop Nests . . . . .	28
6.4	Solution with Multiple IIs . . . . .	28
<b>7</b>	<b>Experiments</b>	<b>30</b>
7.1	Performance of Kernel Loops . . . . .	32
7.1.1	Matrix Multiply . . . . .	33
7.1.2	Tiled and Unroll-and-Jammed Matrix Multiply . . . . .	34
7.1.3	Modified 2-D Explicit Hydrodynamics . . . . .	35
7.2	Performance of SPEC Loops . . . . .	36
7.3	Performance Analysis from Cache Misses . . . . .	37
<b>8</b>	<b>Related Work</b>	<b>38</b>
<b>9</b>	<b>Conclusions</b>	<b>40</b>
<b>A</b>	<b>Appendix: Correctness, Scheduling Algorithm, and Loop Rewriting</b>	<b>44</b>
A.1	Correctness of the Final Schedule with a Single II . . . . .	44
A.2	Correctness of the Final Schedule with Multiple IIs . . . . .	51
A.3	A Scheduling Algorithm . . . . .	53
A.4	Loop Rewriting . . . . .	58

## List of Figures

1	Modulo Scheduling of a Single Loop . . . . .	5
2	A Loop Nest and its DDG . . . . .	6
3	A Conceptual Illustration of Our Software Pipelining Approach (Assume $N_1 = 6$ and $N_2 = 3$ ) . . . . .	6
4	Repeating Patterns in the Final Schedule . . . . .	8
5	Rewritten Loops . . . . .	8
6	Dependences . . . . .	11
7	Dependence Simplification and 1-D Schedule Construction . . . . .	13
8	An Imperfect Loop Nest Example . . . . .	23
9	The Schedules before and after Pushing down (Assume $N_1 = 6$ and $N_2 = 3$ ) . . . . .	24
10	The Rewritten Loop Nest Representing the Final Schedule . . . . .	25
11	The Imperfect Loop Nest Model . . . . .	26
12	The 1-D Schedule with a Single II . . . . .	26
13	The General Form of the 1-D Schedule with Multiple IIs . . . . .	29
14	A 1-D Schedule with 2 IIs for the Loop Nest in Fig. 8(a). There are 2 null cycles above $K_n$ . . . . .	29
15	Compile Flow . . . . .	31
16	Performance of Matrix Multiply, where $N * N$ is the Size of an Array . . . . .	34
17	Tiled MM . . . . .	35
18	Unroll-And-Jammed MM . . . . .	35
19	Hydrodynamics . . . . .	36
20	Percentages of the Three Level Cache Misses in MM- $ikj$ with Matrix Sizes 1024*1024 . . . . .	37
21	L2 Cache Misses . . . . .	38
22	L3 Cache Misses . . . . .	38
23	The Kernel Achieved by the Scheduling Algorithm . . . . .	55
24	Code Generation for An $n$ -dimensional Loop Nest . . . . .	59
25	Illustrating the Code Skeleton. The initialization and finalization phases are not shown, as they are machine-dependent. . . . .	60

# 1 Introduction

Loop nests are rich in coarse-grain and fine-grain parallelism and substantial progress has been made in exploiting the former [4, 8, 10, 18]. With the advent of ILP (Instruction-Level Parallelism) architectures like Very-Long Instruction Word and superscalar processors, and the fast growth in hardware resources, it has been another important challenge to exploit fine-grain parallelism in the loop nests as well.

Software pipelining is an effective way to extract ILP from loops. While numerous algorithms have been proposed for single loops or the innermost loops of loop nests [2, 14, 15, 19, 25, 24], only a few address software pipelining of loop nests [17, 20, 27, 11, 23].

In [17], a loop is modulo scheduled and is considered as an *atomic operation* of its outer loop. The outer loop can then be modulo scheduled. The process is repeated until all loop levels are scheduled, or available resources are used up, or dependences disallow further parallelization. The inefficiency due to the filling and draining (prolog and epilog) of the software pipeline is addressed in [20, 27].

We refer to the above approach as *innermost-loop-centric modulo scheduling*. This approach naturally extends the single loop scheduling method to the multi-dimensional domain, but has two major shortcomings: (1) it commits itself to the innermost loop first without considering how much parallelism the other levels have to offer. Software pipelining another loop level might result in higher parallelism. (2) It cannot exploit the data reuse potential in the outer loops.

There are other software pipelining approaches, developed from hyperplane scheduling. They exploit parallelism from the multi-dimensional iteration space, based on dependences, but cannot handle resource constraints [11, 23].

There has also been other interesting work that combines loop transformations with software pipelining [29, 5]. However, in these methods, software pipelining is still limited to the innermost loop of the transformed loop nest.

This paper presents a framework for resource-constrained software pipelining for a class of loop nests. Software pipelining is applied to the most “beneficial” level in a loop nest, in order to better exploit parallelism and data reuse potential, and match the hardware resources.

The problem addressed in this paper can be formally stated as follows: *Given a loop nest composed of  $n$  loops  $L_1, L_2, \dots, L_n$ , from the outermost to the innermost level, with one loop at each level, identify the most profitable loop  $L_x$  ( $1 \leq x \leq n$ ) that has a rectangular iteration space and software pipeline it.* Software pipelining  $L_x$  means that the consecutive iterations of  $L_x$  will be overlapped at run-time. In this paper, we only discuss how to parallelize the selected loop  $L_x$ . Its outer loops, if any, remain intact in our approach. Since there is only one loop at each level, in this paper, the terms “loop” and “loop level” can be used interchangeably.

The above problem can be broken down into two sub-problems: how to predict the benefits of software pipelining a loop level, and how to software pipeline the most “beneficial” one predicted.

Our solution consists of three steps:

1. *Loop selection*: This step searches for the most profitable loop level in the loop nest. Profitability

can be measured in terms of initiation rate, data reuse potential, or any other optimization criteria. The selected loop may be a loop nest itself, i.e., it may have its own inner loops.

2. *1-D schedule construction*: The multi-dimensional DDG of the selected loop is reduced to a 1-dimensional (1-D) DDG. Based on the 1-D DDG and the resource constraints, a modulo schedule, referred to as a *1-D schedule*, is constructed for the operations in the selected loop. No matter how many inner loops the selected loop has, it is scheduled as if it were a single loop.
3. *Final schedule computation*: Based on the resulting 1-D schedule, this step derives a simple mapping function that specifies the schedule time of the operation instances in the selected loop.

Since the problem of multi-dimensional scheduling is reduced to 1-dimensional scheduling and mapping, we refer to our approach as Single-dimension Software Pipelining (SSP). This approach shows several advantages:

- *Global foresight*: Instead of focusing only on the innermost loop, every loop level is examined and the most profitable one is chosen. Any criterion can be used to judge the “profitability” in this step. This flexibility opens a new prospect to combine software pipelining with any other optimal criterion beyond the ILP degree, which is often the major objective of software pipelining. In this paper, we consider not only parallelism, but also cache effects, which have not been considered by most traditional software pipelining methods.
- *Simplicity*: The method retains the simplicity of the classical modulo scheduling of single loops. The scheduling is based on a simplified 1-dimensional DDG, no matter how many inner loops the selected loop has. This is an essential difference from previous approaches. The traditional modulo scheduling of single loops is subsumed as a special case.
- *Efficiency*: Our schedule divides the iterations of the chosen loop into groups, and executes them group by group. Each group is pipelined. The draining of a group is naturally overlapped with the filling of the next group. For this reason, the schedule is proved to have the shortest length, comparing with any schedule produced by the innermost-centric modulo scheduling under identical conditions for a perfect loop nest. In addition, we search the entire loop nest and choose the most profitable loop level, and consider data reuse in the cache, which also improves the actual execution time of the schedule.
- *Resource sensitivity*: Although our method has been developed independently, we can relate our approach as an extension of hyperplane scheduling [8, 18] in the context of software pipelining for uniprocessors. Our approach extends the traditional hyperplane scheduling to be resource-constrained. The major differences between the two approaches are: (i) Software pipelining is used for uniprocessors, which has limited resources, and it is imperative to consider such constraints. Hyperplane scheduling, however, is usually used for large array-like hardware structures such as systolic and SIMD arrays, and does not consider resource constraints. (ii) Hyperplane scheduling often exploits parallelism from more than one loop level, whereas software pipelining focuses on a single loop level only.

- *Reducing overhead:* Being able to schedule an arbitrary loop level provides freedom for choosing good schedules with less overhead.

To understand this, let us ask a question: why is it necessary to develop a technique that can pipeline a multi-dimensional loop directly? Can't we achieve the same effect through loop transformations followed by innermost-loop-centric software pipelining?

A key insight is that when software pipelining is applied to a loop, the schedule has associated overhead, including initialization, prolog, epilog and finalization. Take the IA-64 architecture as an example. Initialization sets up the initial values of the control registers, including the loop counter register `LC`, epilog counter register `EC`, and predicate rotating register (`pr.rot`). It also transfers live-in values of the loop variables to rotating registers. Finalization transfers live-out values out of rotating registers. Such overhead is incurred every time the loop is executed. The overhead is unavoidable, no matter whatever loop transformations have been performed previously. Some loop transformations, like tiling, would magnify the overhead of an inner loop level, by making the loop nest deeper and the inner loops having smaller trip counts.

Intuitively, the outer the loop is, the less overhead it has. In a 3-deep loop nest, where each loop level is executed 1000 times, the overhead is incurred 1,000,000 times if the innermost loop is software pipelined, 1,000 times if the middle loop is pipelined, and only 1 time if the outermost loop is pipelined. Unless software pipelining the innermost loop results in significant benefit that outweighs such overhead, innermost-loop-centric software pipelining may not be advantageous.

In terms of loop transformations, they are orthogonal to our approach. In this paper, we assume that a loop nest to be software pipelined has already been optimized by loop transformations, if any.

SSP can be applied to both perfect and imperfect loop nests. We will restrict it to perfect loop nests first, and then extend it to imperfect ones, where we introduce multiple initiation intervals into the schedule to improve execution efficiency.

In this paper, we focus only on the fundamental theory of SSP scheduling. The other equally important problems are how to design efficient scheduling algorithms, and how to allocate registers and generate compact code for the constructed SSP schedule. The scheduling algorithm used in our experiments is described in the appendix. Register allocation and code generation have posed interesting challenges and are addressed elsewhere [Rong et al. 2005; Rong and Douillet et al. 2004].

We target ILP uniprocessors with support for predication [15, 3]. Our approach does not impose any constraint on the function units; they may be homogeneous or heterogeneous, pipelined or non-pipelined, and can have unit or multi-cycle latencies.

We have implemented a prototype of our method in the ORC compiler for the IA-64 architecture and tested it with loop nests from Livermore and SPEC2000 floating-point benchmarks. The resulting code is run on an IA-64 Itanium machine and the actual execution time is measured. Preliminary experimental results on these loop nests reveal considerable performance differences in software pipelining different loop levels of a loop nest, and often, software pipelining an outer loop shows significant performance improvements over modulo scheduling that software pipelines the innermost loop. Furthermore, we



observe that SSP is beneficial in the presence of loop transformations, such as loop interchange, loop tiling, and unroll-and-jam.

This paper is organized as follows. Section 2 introduces the basic concepts and briefly reviews modulo scheduling. Then we motivate our study by a simple example in Section 3. Section 4 discusses our method in detail. We prove its correctness and efficiency in Section 5. We then extend SSP to imperfect loop nests and introduce multiple initiation intervals into the schedule in Section 6. Experimental results are reported in Section 7. A discussion on related work and concluding remarks are presented in Sections 8 and 9.

## 2 Basic Concepts

An *n-deep perfect loop nest* is composed of loops  $L_1, L_2, \dots, L_n$ , respectively, from the outermost to the innermost level, with each level having exactly one loop, and all operations are within the innermost loop. Each loop  $L_x (1 \leq x \leq n)$  has an index variable  $i_x$  and a trip count  $N_x > 1$ . The index is normalized to change from 0 to  $N_x - 1$  with unit step. The body of the loop nest consists of all the operations. It is assumed to have no branches; branches, if any, have been converted to linear code by if-conversion [3].

Since a loop, except the innermost loop, has its own inner loops, it is a loop nest itself. To emphasize this fact, we also refer to a loop as an *x-dimensional loop*, where  $x$  is the depth of the loop nest. For example,  $L_1$  is an  $n$ -dimensional loop,  $L_n$  a 1-dimensional loop, etc. A 1-dimensional loop is also called a *single loop*.

The loop nest has an *iteration space*, which contains one point for each execution of the body of the loop nest. Such a point is called an *iteration point* in this paper, and is identified by the index vector  $\mathbf{I} = (i_1, i_2, \dots, i_n)$ . The instance of any operation  $o$  in this iteration point is denoted by  $o(\mathbf{I})$ . The iteration space is *rectangular* if its bounds,  $N_1, N_2, \dots$ , and  $N_n$ , do not change during the execution of the loop nest, although they can change before and after it. In this paper, the loop level to be software pipelined must have a rectangular iteration space.

An  $L_x$  *iteration* is one execution of the  $L_x$  loop body. Thus the  $L_x$  loop has a total of  $N_x$  number of iterations. One such iteration is also an iteration point if  $L_x$  is the innermost loop, i.e.,  $x = n$ .

We use  $(o_1 \rightarrow o_2, \delta, \mathbf{d})$  to represent a data dependence from operation  $o_1$  to operation  $o_2$  in the loop nest, where  $o_1$  and  $o_2$  are called *source* and *sink* of the dependence, respectively;  $\delta \geq 0$  is the *dependence latency*; and  $\mathbf{d} = \langle d_1, d_2, \dots, d_n \rangle$  is the *distance vector*, where  $d_1$  is the distance at the outermost level, and  $d_n$  the innermost.

The *sign of a vector* is that of its first non-zero element, either positive or negative. If all elements are 0, the vector is a *zero vector*.

*Software pipelining* exposes instruction-level parallelism by overlapping successive iterations of a loop. *Modulo scheduling (MS)* is an important and probably the most commonly used approach of software pipelining [14, 17, 25, 24]. A detailed introduction can be found in [2].

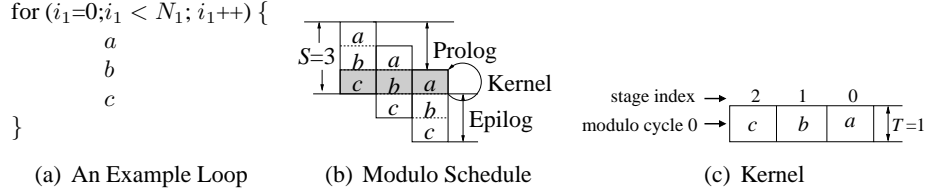


Figure 1: Modulo Scheduling of a Single Loop

Modulo scheduling is usually applied to a single loop. Instances of an operation from successive iterations of the loop are scheduled with an *Initiation Interval* (II) of  $T$  cycles. This is referred to as *modulo property*. A valid modulo schedule respects modulo property, the *dependence constraints*, and the (*hardware*) *resource constraints*.

The schedule length  $l$  is defined as the execution time of a single iteration. Then each iteration is composed of  $S = \lceil \frac{l}{T} \rceil$  number of *stages*, with each stage taking  $T$  cycles. The schedule consists of three phases: the *prolog* to fill the pipeline, the *kernel* to be executed multiple times, and the *epilog* to drain the pipeline.

Let the schedule time for any operation instance  $o(i_1)$  be  $\sigma(o, i_1)$ . When  $i_1 = 0$ , it can be expressed as  $\sigma(o, i_1) = p * T + q$ , where  $0 \leq q < T$ . We say that operation  $o$  is scheduled to *modulo cycle*  $q$  within stage  $p$ .

**Example:** Fig. 1(a) shows an example loop. Assume 3 function units and two dependences ( $a \rightarrow b, 1, \langle 0 \rangle$ ) and ( $b \rightarrow c, 1, \langle 0 \rangle$ ). Fig. 1(b) shows a modulo schedule for the loop with  $T = 1$ , and  $S = 3$ . Fig. 1(c) specifically shows the kernel, where the stages are numbered from 0 to 2 from right to left, and all operations are scheduled to modulo cycle 0.

### 3 Motivation and Overview

In this section, we motivate our method with the help of a simple 2-deep perfect loop nest. We bring out the practical limitations of the innermost-centric approach and motivate the necessity of our work. Subsequently, we illustrate our approach using this example. Then we summarize the intuitions we get from the example, and briefly describe our theoretical solution to the general problem.

#### 3.1 A Motivating Example

Fig. 2 shows a perfect loop nest in C language and its data dependence graph. This loop nest certainly could be parallelized in a number of other ways, too. We use it only for illustration purposes.

To facilitate understanding and without loss of generality, in this example, we assume that each statement is an operation. In the DDG, each node represents an operation and an edge represents a dependence labeled with the distance vector.

The inner loop has no parallelism due to the dependence cycle  $a \rightarrow b \rightarrow a$  at this level. Thus modulo scheduling of the inner loop cannot find any parallelism for this example. Innermost-loop-centric

```

L1: for (i1=0; i1 < N1; i1++){
L2:  for (i2=0; i2 < N2; i2++){
      a: U[i1 + 1][i2]=V[i1][i2]+U[i1][i2];
      b: V[i1][i2 + 1]=U[i1 + 1][i2];
    }
  }

```

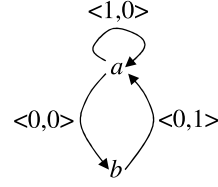


Figure 2: A Loop Nest and its DDG

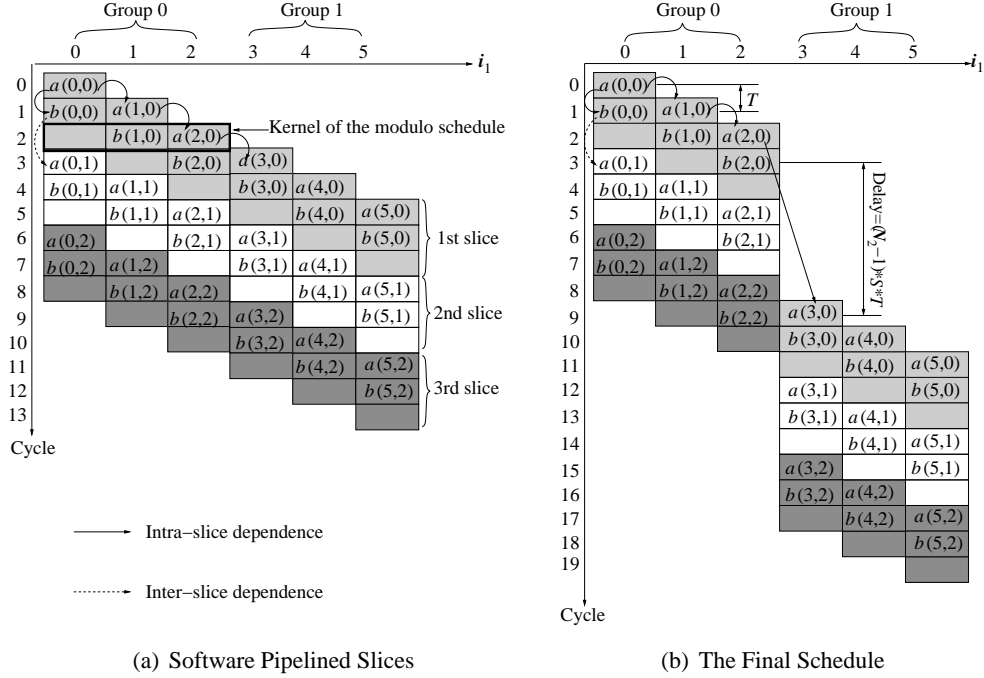


Figure 3: A Conceptual Illustration of Our Software Pipelining Approach (Assume  $N_1 = 6$  and  $N_2 = 3$ )

software pipelining approach exposes extra parallelism by overlapping the filling and draining of the pipeline between successive outer loop iterations. Since modulo scheduling failed to find any parallelism, there is no filling or draining and therefore no overlapping. Thus, innermost-loop-centric software pipelining cannot find any parallelism, either.

One may argue that loop interchange before software pipelining will solve this problem. Unfortunately, that will destroy the data reuse in the original loop nest: for large arrays each iteration point will introduce 2 cache misses, as the array elements are now accessed column-wise rather than row-wise.

### 3.2 Illustration of Our Approach

The above example shows the limitation of the traditional software pipelining: It cannot see the whole loop nest to better exploit parallelism. Nor can it exploit the data reuse potential of the outer loop(s).

This raises the question: Why not select a better loop to software pipeline, not necessarily the innermost one?

This question brings the challenging problem of software pipelining of a loop nest. The challenge comes from two aspects: how to handle resource constraints? And how to handle the multi-dimensional dependences? Before we expand discussion on these challenges, let us once again look at our motivating example shown in Fig. 2.

**Example:** Let us assume operations  $a$  and  $b$  have latencies of 1 and 2 cycles, respectively. Assume that we have two functional units, and both are pipelined and can perform any of the operations.

Suppose that the outer loop  $L_1$  is selected for software pipelining. We remember software pipelining a loop is to overlap its iterations. Fig. 3(a) shows such an overlapping, where the initiation interval between two adjacent iterations of the loop is  $T = 1$  cycle, and we assume  $N_1 = 6$  and  $N_2 = 3$  for simplicity.

We consider the operations belonging to iteration points  $(i_1, 0)$ , for all  $i_1$ , constitute the first *slice*, and operations belonging to points  $(i_1, 1)$  the second slice, etc. Then the overlapping can be reinterpreted in this way: each slice is modulo scheduled so that successive iteration points within this slice initiate at an interval of  $T = 1$  cycle. For the first slice, the kernel of the modulo schedule is highlighted in a box. There are  $S = 3$  stages, with one stage being empty.

Although the resource constraints are respected within each modulo scheduled slice, they are violated between slices because a slice is issued greedily without waiting for the resources to be released by the previous slice. To remove the conflicts, we cut the slices into *groups*, with each group having  $S = 3$  iterations of the outer loop. There are two groups in this schedule. Each group, except the first one, is pushed down by  $(N_2 - 1) * S * T$  cycles relative to its previous group. The delay is designed to ensure that repeating patterns definitely appear. This leads to the *final schedule* that maps each instance of an operation to its schedule time, as shown in Fig. 3(b). Note that not only dependence and resource constraints are respected, but the parallelism degree exploited in a modulo scheduled slice ( $S = 3$ ) is still preserved, and the resources are fully used. A dependence is still respected after the pushing-down because that action either does not affect, or only increases, the time distance between the source and the sink of the dependence, as illustrated by the dependences in Fig. 3(a) before the pushing-down and in Fig. 3(b) after that.

Repeating patterns can be found in the final schedule. In Fig. 4, we add to the final schedule some ineffective operation instances, as shown in the shaded part. They are ineffective because their first indexes are beyond the legal range of  $i_1$ , the outer loop index variable. The range is assumed to be  $[0, 6)$  in our illustration. For target architectures with predication support like IA-64, predicate registers can be used to make them ineffective during execution of the final schedule [Rong and Douillet et al. 2004]. With the added ineffective operation instances, it is clear to see that the final schedule is composed of two repeating patterns, referred to as *Outermost Loop Pattern (OLP)* and *Inner Loop Execution Segment (ILES)*. An OLP drains the pipeline of a group, and fills the pipeline with the next group simultaneously. When the pipeline is filled with the next group, an ILES starts. It runs all the inner loops of the next group until the group is going to drain. Then another OLP starts. Note that an ILES itself is composed of  $N_2 - 1 = 2$  number of a smaller pattern, as shown in the figure. Apart from the OLPs and ILESes,



### 3.3 Overview of Our Approach

Based on the illustration, in this section, we briefly describe the steps and challenges in solving our general problem. The first challenge is how to select a loop level for software pipelining. Once the loop level is identified, the second challenge is how we software pipeline it, taking into account resource and dependence constraints. The principles are discussed below, while details are left to Section 4.

#### 3.3.1 Which Loop to Software Pipeline?

Parallelism is surely one of the major concerns. On the other hand, cache effects are also important and govern the actual execution time of the schedule. However, it is hard to consider cache effects in traditional software pipelining, mainly due to the fact that it focuses on the innermost loop. Provided that an arbitrary loop level in a loop nest can be software pipelined, we can search for the most profitable level, measured by parallelism or cache data reuse, or both. Any other objective can also be used as a criterion.

The selected loop, which is a loop nest itself, needs to have a rectangular iteration space. How to handle a loop with a non-rectangular iteration space is beyond the scope of this paper.

#### 3.3.2 How to Software Pipeline the Selected Loop?

Suppose we have chosen a loop, for simplicity, say, the outermost loop  $L_1$ . Conceptually, we allocate the iteration points within the loop to a series of *slices*, and software pipeline each slice. Although any software pipelining method can be used, we focus on modulo scheduling in this paper.

The iteration points are allocated in this way: for any  $i_1 \in [0, N_1)$ , iteration point  $(i_1, 0, \dots, 0, 0)$  is allocated to the first slice,  $(i_1, 0, \dots, 0, 1)$  to the second slice, and so on.

To modulo schedule a slice, we reduce the DDG to have only the dependences at this  $L_1$  loop level and simplify their distance vectors to be 1-dimensional. Based on the resource constraints and this 1-D DDG, we construct a modulo schedule, referred to as a *1-D schedule*.

This 1-D schedule is applied to every slice. So all slices have the same shape, and they can be packed together seamlessly. This leads to a schedule like that in Fig. 3(a). We can see that the iteration points are allocated to the slices in such a way that all iterations of the  $L_1$  loop run in parallel, while each of them runs sequentially.

*How to Handle Resources?* Resource constraints are enforced at two levels: at the slice level when we modulo schedule the slices, and at the inter-slice level. Modulo scheduling of a slice meets the resource constraints within the slice. However, by packing the successive slices together, two slices are partially overlapped. The collective resources required at an overlapping cycle may exceed the resources available. To solve this problem, we cut the slices into *groups*, with each group containing  $S$  iterations of the  $L_1$  loop. Then we push down, i.e. delay the execution of, a group until resources are available. This results in a *final schedule*, which respects the resource constraints at each cycle, as illustrated in Fig. 3(b).

*How to Handle Dependences?* A major obstacle to software pipelining of a loop nest is how to handle the  $n$ -dimensional dependence distance vectors. As mentioned earlier, in modulo scheduling a slice, we consider only the simplified dependences with 1-dimensional distance vectors. Doing so is sufficient to satisfy all dependence constraints within a slice. Dependences across slices (in the forward direction) are also satisfied, since the slices are executed sequentially. After the slices are cut into groups and the groups are pushed down, the dependences are still respected, as pushing down either does not affect, or can only increase, the time distance between the source and the sink of a dependence, as illustrated in Fig. 3(a) and Fig. 3(b).

### 3.3.3 Constructing the Final Schedule

In this paper, we express the final schedule abstractly in a function, which is *independent* of any specific architecture. The function describes the final schedule time of an operation instance, based on the 1-D schedule time of that operation.

It is important to see that this function is determined by the 1-D schedule. We do *not* unroll any loop in constructing either the 1-D schedule or the final schedule. The example in Fig. 3(a) and 3(b) has illustrated the formation of the final schedule in a way that can be easily understood: fully unrolling the chosen loop and allocating all the iteration points in it to slices, applying the 1-D schedule to all slices, cutting them into groups, and pushing down the groups appropriately. In our formal solution, however, the same effect is simply captured by the final schedule function.

For a specific architecture, the final schedule is constructed by composing the prolog, OLP, ILES, and epilog with the 1-D schedule. This realizes the function equivalently by code. Depending on the target architecture, the 1-D schedule may need to be duplicated in this process. For example, for an architecture like IA-64 that supports rotating register files, the final schedule rotates registers in an OLP, but stops rotating in an ILES. Thus the ILES has to duplicate the 1-D schedule to achieve the same effect of register rotating. The details of this code generation process are beyond the scope of this paper and are presented elsewhere [Rong and Douillet et al. 2004].

Later, when we extend our approach to allow all the loops in the loop nest have their own distinct IIs, it is too complicated to express the final schedule in a function. In this case, we also resort to the constructive code generation process.

In summary, our approach to software pipeline a loop nest consists of three steps: loop selection, 1-D schedule construction, and final schedule computation. We will describe them in detail in the next section.

## 4 Solution

In this section, we first define the concept of simplified DDG. With this concept, we formalize our approach into 3 steps: (1) loop selection, (2) 1-D schedule construction<sup>1</sup>, and (3) final schedule com-

---

<sup>1</sup>In our previous work [Rong and Tang et al. 2004], this step is referred to as “dependence simplification and 1-D schedule construction”. We remove “dependence simplification” from the name, because dependence simplification is also needed by

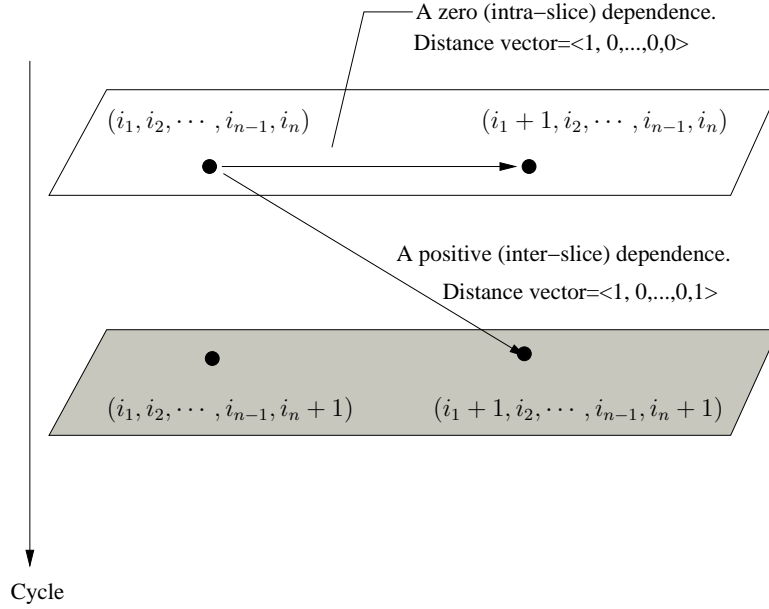


Figure 6: Dependences

putation.

#### 4.1 Definition of Simplified DDG

As illustrated in Fig. 3, conceptually, the final schedule of a multi-dimensional loop consists of a series of modulo scheduled slices, which are cut into groups, and the groups are then pushed down to resolve inter-slice resource conflicts. If a dependence is respected before pushing down the groups, it will also be respected after that. Therefore we only need to consider the dependences necessary to obtain the modulo schedule before the pushing-down.

Fig.6 pictorially illustrates the dependences in an  $n$ -dimensional loop nest in two successive slices, where each parallelogram represents a slice, and each dot an iteration point. Although not shown on the picture, each slice is software pipelined. The outermost level  $L_1$  is assumed to be the chosen loop. There are two kinds of dependences: one is across two slices, and the other one is within a slice.

Due to the way the iteration points are allocated, a dependence across two slices has a distance vector  $\langle d_1, d_2, \dots, d_n \rangle$ , where  $d_1 \geq 0$ , and  $\langle d_2, \dots, d_n \rangle$  is a positive vector. Such a dependence is naturally resolved because the two slices are executed sequentially.

A dependence within a slice has a distance vector  $\langle d_1, d_2, \dots, d_n \rangle$ , where  $d_1 \geq 0$ , and  $\langle d_2, \dots, d_n \rangle$  is a zero vector. Such a dependence has to be considered during software pipelining of the slice. Besides, only the distance  $d_1$  is useful for software pipelining. That is, the dependence distance vector can be simplified as  $\langle d_1 \rangle$  in pipelining.

---

the loop selection step. Therefore, it is more appropriate to be taken as a basic concept shared by both steps.



The two kinds of dependences are named *positive* and *zero dependences*, respectively. Note that a dependence from a slice to a previous slice is illegal. It is called a *negative dependence*. Negative dependences cannot be handled directly<sup>2</sup>.

Below we formally classify the dependences, and define the simplified dependence graph.

Let  $\mathbf{d} = \langle d_1, d_2, \dots, d_n \rangle$  be the distance vector of a dependence. We say that this dependence is *effective at loop level*  $L_x (1 \leq x \leq n)$  iff  $\langle d_1, d_2, \dots, d_{x-1} \rangle = \mathbf{0}$  and  $\langle d_x, d_{x+1}, \dots, d_n \rangle \geq \mathbf{0}$ , where  $\mathbf{0}$  is the zero vector with appropriate length. By *effective*, we mean that such a dependence must be respected by the final schedule if we software pipeline  $L_x$ . All effective dependences at  $L_x$  compose the *effective DDG* at  $L_x$ .

According to the definition, if a dependence is effective at  $L_x$ , we have  $\langle d_x, d_{x+1}, \dots, d_n \rangle \geq \mathbf{0}$ . Of course the first element  $d_x \geq 0$ . We classify the dependence by the sign of the sub-distance-vector  $\langle d_{x+1}, \dots, d_n \rangle$ , when  $x < n$ . If this sub-vector is a zero, positive, or negative vector, the dependence is classified as a *zero, positive, or negative dependence at  $L_x$* , respectively. When  $x = n$ , we classify it as a zero dependence at  $L_x$  for uniformity.

The above classification is complete: *an effective dependence is in and only in one of the three classes*. Especially, the dependences are classified according to *the sign of the sub-distance-vector*, not that of the whole distance vector. For example, a dependence in a 3-deep loop nest with a distance vector of  $\langle 1, -1, 2 \rangle$  is a negative dependence at  $L_1$  because the sub-vector  $\langle -1, 2 \rangle$  is negative, even though the whole distance vector is positive.

We classify only effective dependences, since, in the following sections, our discussion relates only to them. Although the dependence classification is dependent on the loop level, we will not mention the loop level when the context is clear.

In this paper, we assume that when we consider to software pipeline a loop level  $L_x$ , all effective dependences at this level are either zero or positive. As discussed above, a positive dependence is across slices in the forward direction and can be naturally honored in the final schedule. Only zero dependences are within a slice and need to be considered. Lastly, only the dependence distance at  $L_x$  is useful for software pipelining. Thus we can reduce the effective DDG to have only zero dependences with 1-dimensional distance vectors. We refer to the resulting DDG as the *simplified DDG*. The definition is as follows: The *simplified DDG* at  $L_x$  is composed of all the zero dependences at  $L_x$ ; the dependence arcs are annotated with the dependence distance at  $L_x$ .

**Example:** Fig.7(a) shows the effective DDG at  $L_1$  for the loop nest depicted in Fig.2. There are two zero dependences in this DDG:  $a \rightarrow a$  and  $a \rightarrow b$ . Associating the dependence distances at  $L_1$  with the arcs, we get the simplified DDG shown in Fig.7(b).

---

<sup>2</sup>The loop nest must be transformed to make negative dependences become zero or positive. It is always feasible to do so by loop skewing. However, after that, the iteration space becomes non-rectangular. Although we restrict to rectangular iteration spaces in this paper, the first two steps of SSP are still applicable to non-rectangular cases, without any change, since scheduling considers only DDG and hardware resources. It considers nothing about the shape of the iteration space. For the third step of SSP, predicate registers can be used to dynamically form a non-rectangular iteration space in runtime. Another way to handle the non-rectangular iteration space is to apply loop peeling such that the space is cut into a rectangle with a triangle before and after it, and SSP is applied only to the rectangle.

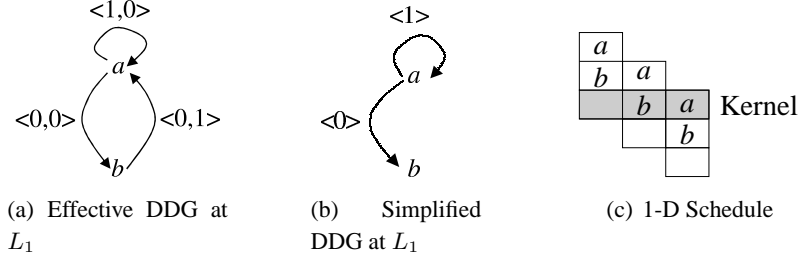


Figure 7: Dependence Simplification and 1-D Schedule Construction

## 4.2 Step 1: Loop Selection

In this paper, our objective is to generate the most efficient software-pipelined schedule for a loop nest. Thus it is desirable to select the loop level with a higher initiation rate (higher parallelism), or a better data reuse potential (better cache effect), or both. The specific decision is not made here, since that is implementation-dependent. This paper focus on presenting a general framework, not algorithm or implementation. In this section, we address the essential problem of evaluating these two criteria. For each criterion, we consider all the loop levels that have rectangular iteration spaces.

### 4.2.1 Initiation Rate

*Initiation rate*, which is the inverse of initiation interval, specifies the number of iteration points issued per cycle. Hence we choose the loop level  $L_x$  that has the maximum initiation rate, or minimum initiation interval.

The minimum initiation interval at loop level  $L_x$  is  $\max(RecMII_x, ResMII)$ , where  $RecMII_x$  and  $ResMII$  are the minimum initiation intervals determined, respectively, by recurrences in the simplified DDG at  $L_x$ , and by the available hardware resources<sup>3</sup>.

$$RecMII_x = \max_{\forall C} \frac{\delta(C)}{d(C)}, \quad (1)$$

where  $C$  is a cycle in the simplified DDG,  $\delta(C)$  is the sum of the dependence latencies along cycle  $C$ , and  $d(C)$  is the sum of the dependence distances along  $C$  [13].

$$ResMII = \max_{\forall \text{ resource type } r} ResMII_r \quad (2)$$

where  $ResMII_r$  is the lower bound of MII determined by resource type  $r$ , which is

<sup>3</sup>The reader will find in section 4.3 that our approach has extra Sequential Constraints. They affect only the schedule length of the 1-D schedule, but not the initiation interval. In the worst case, we can always increase the schedule length to satisfy these constraints. Thus they do not influence the MII calculation here.

$$ResMII_r = \begin{cases} \frac{\text{total operations that use } r}{\text{total resources of type } r} & \text{if } r \text{ is pipelined.} \\ \frac{\text{total execution time of the operations that use } r}{\text{total resources of type } r} & \text{if } r \text{ is non-pipelined.} \end{cases} \quad (3)$$

In addition to the initiation rate, we also look at the trip count of each loop level. In particular, the trip count should not be less than  $S$ , the number of stages in the 1-D schedule. Otherwise, this loop should not be chosen.

The reason is that the slices are cut in groups, where each group has  $S$  iterations of loop  $L_x$ . Then the trip count  $N_x$  is expected to be divisible by  $S$ . Otherwise, the last group will have fewer  $L_x$  iterations, resulting in a lower utilization of resources in that group. However, when  $N_x > S$ , it is always possible to apply loop peeling to avoid the situation.

Although  $S$  is unknown at loop selection time, it is generally small because the limited resources in a uniprocessor cannot support too many stages at the same time. As a guideline, a small estimated value can be set for  $S$ .

#### 4.2.2 Data Reuse

When we software pipeline a loop level, the data reuse potential can be measured by the average number of memory accesses per iteration point. The fewer the accesses, the greater the reuse potential. Without loss of generality, let us consider loop  $L_1$ .

In our approach, software pipelining results in  $S$  iterations of  $L_1$  loop running in a group, which is composed of a series of slices. Select the first  $S$  number of successive slices in the first group. They include the following set of iteration points:  $\{(i_1, 0, \dots, 0, i_n) | \forall i_1 \text{ and } i_n \in [0, S)\}$ , which is an  $S * S$  square in the iteration space. This is a typical situation in our method, because  $L_1$  iterations are executed in parallel, and the index of the innermost loop changes more frequently than the indices of the other loops. Therefore, we could estimate the memory accesses of the whole loop nest by those of the iteration points in this set. This set can be abstracted as a *localized vector space*  $\alpha = span\{(1, 0, \dots, 0), (0, \dots, 0, 1)\}$ . Now the problem is very similar to that discussed in [28]. Below we briefly describe the application of their method in this situation.

For a *uniformly generated set* of memory references in this localized space, let  $R_{ST}$  and  $R_{SS}$  be the self-temporal and self-spatial reuse vector spaces, respectively. And let  $gT$  and  $gS$  be the number of group temporal and group-spatial equivalent classes. Then for this uniformly generated set, the number of memory accesses per iteration point is <sup>4</sup>:

$$\frac{gS + (gT - gS)/l}{l^e \mathcal{S}^{\dim(R_{ST} \cap \alpha)}}, \quad (4)$$

where  $l$  is the cache line size, and

---

<sup>4</sup>We use  $R_{ST}$  in the denominator here, instead of  $R_{SS}$  as in the original formula on page 39 of the literature[28], which we think was a typo.

$$e = \begin{cases} 0 & \text{if } R_{ST} \cap \alpha = R_{SS} \cap \alpha, \\ 1 & \text{otherwise.} \end{cases}$$

The total number of memory accesses per iteration point is the sum of accesses for each uniformly generated set.

The above data reuse model does not consider loop volume. Other models [12, 7, 16] may be used as well.

### 4.3 Step 2: 1-D Schedule Construction

Our method software pipelines only the selected loop. Enclosing outer loops, if any, are left as they are. Therefore, without loss of generality, we consider  $L_1$  as the selected loop.

As mentioned already, given the effective DDG at  $L_1$ , we can simplify the dependences to obtain a simplified DDG, which consists of only zero dependences with 1-dimensional distance vectors.

Based *solely* on the simplified DDG and the hardware resource constraints, we construct a 1-D schedule. Since the DDG is 1-dimensional, from the viewpoint of scheduling,  $L_1$  is treated as if it were a single loop. Any modulo scheduling method can be applied to obtain the 1-D schedule.

Let  $T$  be the initiation interval of the generated schedule, and  $S$  be the number of stages of the schedule. We refer to the schedule as a *1-D schedule* for the loop level  $L_1$ . Let the schedule time for any operation instance  $o(i_1)$  be  $\sigma(o, i_1)$ , where  $0 \leq \sigma(o, i_1) < S * T$  when  $i_1 = 0$ .

The 1-D schedule must satisfy the following properties:

1. Modulo property:

$$\sigma(o, i_1) + T = \sigma(o, i_1 + 1) \quad (5)$$

2. Dependence constraints:

$$\sigma(o_1, i_1) + \delta \leq \sigma(o_2, i_1 + k) \quad (6)$$

for every dependence  $(o_1 \rightarrow o_2, \delta, \langle k \rangle)$  in the simplified DDG.

3. Resource constraints: At any modulo cycle of the kernel, no hardware resource is allocated to more than one operation.
4. Sequential Constraints: if  $n > 1$ , then

$$S * T - \sigma(o, 0) \geq \delta \quad (7)$$

for every positive dependence with operation  $o$  as the source operation, and  $\delta$  being the dependence latency.

The first three constraints are exactly the same as those of the classical modulo scheduling [2, 13]. We have added the sequential constraints to enforce sequential order between successive iteration points in the same  $L_1$  iteration. This ensures that all positive dependences are honored at runtime. The sequential constraints have effect only for loop nests with more than 1 loop.

**Example:** For the loop nest in Fig.2 and its effective DDG in Fig.7(a), the simplified DDG at  $L_1$  is shown in Fig.7(b). Based on this simplified DDG, a 1-D schedule can be constructed (Fig. 7(c)). As mentioned earlier, we have assumed two homogeneous functional units, and an execution latency of 1 and 2 cycles for operations  $a$  and  $b$ . The schedule has an initiation interval of 1 cycle ( $T = 1$ ) and has 3 stages ( $S = 3$ ). Also,  $\sigma(a, i_1) = 0 + i_1 * T$  and  $\sigma(b, i_1) = 1 + i_1 * T$ .

#### 4.4 Step 3: Final Schedule Computation

As explained in Section 3.3.2, we first allocate iteration points in the loop nest to slices, then we software pipeline each slice by applying the 1-D schedule to it.

If the successive slices are greedily issued without considering resource constraints across the slices, we obtain the schedule like that in Fig.3(a). Note that, within each slice, the resource constraints are honored during the construction of the 1-D schedule. Now, to enforce resource constraints across slices, we cut the slices in groups, with each group having  $S$  number of  $L_1$  iterations. Each group, except the first one, is delayed by a given number of cycles as shown in Fig.3(b).

With the above procedure in mind, a final schedule can be precisely defined by the following mapping function. For any operation  $o$  in the iteration point  $\mathbf{I}=(i_1, i_2, \dots, i_n)$ , the schedule time  $f(o, \mathbf{I})$  is given by

$$\begin{aligned}
f(o, \mathbf{I}) &= \sigma(o, i_1) \\
&+ \sum_{2 \leq x \leq n} (i_x * (\prod_{x < y \leq n+1} N_y) * S * T) \\
&+ \left\lfloor \frac{i_1}{S} \right\rfloor * ((\prod_{2 \leq x \leq n+1} N_x) - 1) * S * T, \tag{8}
\end{aligned}$$

where  $N_{n+1}=1$ .

Let us briefly explain how the above equation is derived. First, let us consider the ideal schedule before pushing down the groups. For this schedule, the schedule time of  $o(\mathbf{I})$  is equal to that of  $o(i_1, 0, \dots, 0)$  plus the time elapsed between the schedule times of  $o(i_1, 0, \dots, 0)$  and  $o(\mathbf{I})$ . Since  $o(i_1, 0, \dots, 0)$  is in the first slice, the schedule time of  $o(i_1, 0, \dots, 0)$  is simply equal to  $\sigma(o, i_1)$ , the 1-D schedule time. This corresponds to the first term of the right-hand side of Equation (8).

Between iterations  $o(i_1, 0, \dots, 0)$  and  $o(i_1, i_2, \dots, i_n)$ , there are  $i_2*(N_3*N_4*..*N_n) + i_3*(N_4*N_5*..*N_n) + \dots + i_n$  number of iteration points. These points execute sequentially and each of them takes  $S*T$  cycles. Thus, the time elapsed between the schedule times of  $o(i_1, 0, \dots, 0)$  and  $o(i_1, i_2, \dots, i_n)$  equals

$$\sum_{2 \leq x \leq n} (i_x * (\prod_{x < y \leq n+1} N_y) * S * T).$$

This corresponds to the second term of the right-hand side of Equation (8).

Next we discuss the effect of pushing down the groups. Iteration point  $o(\mathbf{I})$  is located in group  $\lfloor \frac{i_1}{S} \rfloor$ . Each group is delayed by  $\lfloor \frac{i_1}{S} \rfloor * w$  cycles, where  $w$  is the delay between two successive groups. For the example in Fig.3(b) with the 2-deep loop nest, we see that  $w = (N_2 - 1) * S * T$ . In general, for an  $n$ -deep loop nest,  $w = (\text{total iteration points in an } L_1 \text{ iteration} - 1) * S * T$ . Thus the group where  $o(\mathbf{I})$  is located is pushed down by

$$\left\lfloor \frac{i_1}{S} \right\rfloor * \left( \prod_{2 \leq x \leq n+1} N_x - 1 \right) * S * T$$

cycles. This is exactly the third term in Equation (8).

**Example.** To illustrate the mapping function for the final schedule, consider the 2-deep loop nest in Fig.2. From the 1-D schedule in Fig.7(c), we know that  $S = 3$ ,  $T = 1$ , and  $\sigma(a, i_1) = 0 + i_1 * T$ . For any operation instance  $a(i_1, i_2)$ , we have the final schedule

$$f(a, (i_1, i_2)) = i_1 + i_2 * 3 + \left\lfloor \frac{i_1}{3} \right\rfloor * (N_2 - 1) * 3.$$

For instance, when  $N_2 = 3$ , we have  $f(a, (4, 1)) = 13$ , as can be seen from Fig.3(b).

## 5 Analysis

In this section, we establish the correctness and efficiency of the SSP approach, and its relationship with MS. We also demonstrate the relationship between SSP and the traditional hyperplane scheduling.

### 5.1 Correctness

First, we show a simple fact in an SSP final schedule: the instances of any operation are initiated one by one every  $T$  cycles. For example, in Fig. 3(b), the instances of operation  $a$  are issued in a sequence every  $T = 1$  cycle:  $a(0, 0)$ ,  $a(1, 0)$ ,  $a(2, 0)$ ,  $a(0, 1)$ ,  $a(1, 1)$ ,  $a(2, 1)$ ,  $a(0, 2)$ ,  $a(1, 2)$ ,  $a(2, 2)$ ,  $a(3, 0)$ ,  $a(4, 0)$ ,  $a(5, 0)$ ,  $a(3, 1)$ , ... It is trivial to prove that such initiation pattern is true in general for any-deep loop nest with any number of operations.

A direct consequence of the fact is that no two instances of the same operation can be initiated at the same cycle. This result will be used in proving the following theorem:

**Theorem 5.1.** *The final schedule defined in Equation (8) respects all the dependences in the effective DDG and the resource constraints.*

*Proof.* Given a dependence  $(a \rightarrow b, \delta, \langle d_1, d_2, \dots, d_n \rangle)$  in the effective DDG,  $\mathbf{I} = (i_1, i_2, \dots, i_n)$ , and

$\mathbf{I}' = (i_1 + d_1, i_2 + d_2, \dots, i_n + d_n)$ , we show that  $f(b, \mathbf{I}') - f(a, \mathbf{I}) \geq \delta$ . Consider

$$\begin{aligned}
f(b, \mathbf{I}') - f(a, \mathbf{I}) &= \underbrace{\sigma(b, i_1 + d_1) - \sigma(a, i_1)}_{(A)} + \\
&\quad \underbrace{\sum_{2 \leq x \leq n} (d_x * (\prod_{x < y \leq n+1} N_y) * S * T)}_{(B)} + \\
&\quad \underbrace{\left( \left\lfloor \frac{i_1 + d_1}{S} \right\rfloor - \left\lfloor \frac{i_1}{S} \right\rfloor \right) * \left( \left( \prod_{2 \leq x \leq n+1} N_x \right) - 1 \right) * S * T}_{(C)}. \tag{9}
\end{aligned}$$

First, if this is a zero dependence, then  $(B) = 0$ , and  $(C) \geq 0$ . Thus, using Inequality (6), we have

$$f(b, \mathbf{I}') - f(a, \mathbf{I}) \geq \sigma(b, i_1 + d_1) - \sigma(a, i_1) \geq \delta.$$

Therefore zero dependences are respected in the final schedule. On the other hand, if the dependence is positive, then it is easy to see that  $(A) = \sigma(b, d_1) - \sigma(a, 0) \geq -\sigma(a, 0)$ ,  $(B) \geq S * T$ , and  $(C) \geq 0$ . So

$$\begin{aligned}
f(b, \mathbf{I}') - f(a, \mathbf{I}) &\geq S * T - \sigma(a, 0) \\
&\geq \delta \text{(by sequential constraints in Inequality 7)}
\end{aligned}$$

Therefore positive dependences are also respected in the final schedule.

Lastly, any two operation instances that have the same final schedule time must come from the same modulo cycle in the kernel, but they cannot be the instances of the same operation. Since the kernel contains exactly one instance for each operation, and is free of resource competition (by the resource constraints definition in Section 4.3), these operation instances have no resource contention either.

To show this, consider two distinct operation instances,  $a(\mathbf{I})$  and  $b(\mathbf{I}')$ , scheduled at the same cycle. Then

$$f(b, \mathbf{I}') - f(a, \mathbf{I}) = 0.$$

By combining this with Equation (9), we get  $(A) + (B) + (C) = 0$ . Since  $(A) = \sigma(b, i_1) + d_1 * T - \sigma(a, i_1)$ , and  $(B) + (C)$  is a multiple of  $T$  (say,  $p * T$ , where  $p$  is an integer), we have  $\sigma(b, i_1) - \sigma(a, i_1) = (-p - d_1) * T$ , which is also a multiple of  $T$ . This means operations  $a$  and  $b$  must be from the same modulo cycle in the kernel.

As discussed above, the instances of the same operation are issued in order. No two of them can have the same schedule time. Therefore,  $a$  and  $b$  must be different operations.  $\square$

## 5.2 Efficiency

Next, we demonstrate the efficiency of the SSP approach over other innermost-loop-centric software pipelining methods from the viewpoint of computation time of the constructed schedule. In particular, we compare our approach with modulo scheduling of the innermost loop (MS), and modulo scheduling

of the innermost loop and overlapping the filling and draining of adjacent iterations of the outer loop, referred to as extended modulo scheduling (xMS) in this paper [17, 20, 27]. Let us define the *computation time* as the (final) schedule time of the last operation instance+1.

**Theorem 5.2.** *For an  $n$ -deep perfect loop nest, suppose that MS, xMS, and SSP find the same initiation interval  $T$  and stage number  $S$ . Furthermore, suppose that SSP chooses the outermost loop  $L_1$ , which has a trip count  $N_1$ . If  $N_1$  is divisible by  $S$ , then the computation time of the SSP final schedule is not longer than that of the MS or xMS schedule.*

*Proof.* Modulo scheduling parallelizes the innermost loop, whose iterations issue once every  $T$  cycles. So the computation time is

$$Time_{MS} = N_1 * \dots * N_{n-1} * (N_n + S - 1) * T. \quad (10)$$

In the best case of xMS schedule, the cost of filling and draining the pipeline is incurred only at the beginning and end of the execution of the entire loop nest, and an iteration point is issued every  $T$  cycles. The computation time is then

$$Time_{xMS} = (N_1 * \dots * N_n + S - 1) * T. \quad (11)$$

Let  $o$  be any operation and  $\mathbf{I}=(i_1, i_2, \dots, i_n)$  be any index vector. In our approach, it is easy to see that  $f(o, \mathbf{I})$  is maximum when  $\mathbf{I}=(N_1-1, N_2-1, \dots, N_n-1)$  and  $\sigma(o, 0) = S * T - 1$ . The computation time of SSP is equal to the maximal  $f(o, \mathbf{I})+1$ , which is

$$\begin{aligned} Time_{SSP} &= (S + N_1 - 1) * T + \\ &\quad \underbrace{\sum_{2 \leq x \leq n} ((N_x - 1) * (\prod_{x < y \leq n+1} N_y) * S * T)}_{(D)} + \\ &\quad \underbrace{\left\lfloor \frac{N_1 - 1}{S} \right\rfloor * ((\prod_{2 \leq x \leq n+1} N_x) - 1) * S * T}_{(E)}. \end{aligned} \quad (12)$$

It is easy to show that

$$(D) = ((\prod_{2 \leq x \leq n+1} N_x) - 1) * S * T. \quad (13)$$

Further, under the given condition that  $N_1$  is divisible by  $S$ , we know that

$$N_1 - 1 = p * S + (S - 1),$$

where  $p$  is an integer. Thus

$$(E) = p * ((\prod_{2 \leq x \leq n+1} N_x) - 1) * S * T. \quad (14)$$



Combining Equations (14) and (13), we get

$$\begin{aligned} (D) + (E) &= (p + 1) * ((\prod_{2 \leq x \leq n+1} N_x) - 1) * S * T \\ &= N_1 * ((\prod_{2 \leq x \leq n+1} N_x) - 1) * T. \end{aligned}$$

Substituting it in Equation (12), we get

$$Time_{SSP} = ((\prod_{1 \leq x \leq n+1} N_x) + S - 1) * T. \quad (15)$$

From Equations (10), (11), and (15), we have:

$$Time_{SSP} = Time_{xMS} \leq Time_{MS}.$$

□

Intuitively, this theorem holds because the final schedule produced by SSP always issues one iteration point every  $T$  cycles, without any hole, as can be seen from the example in Fig 3(b).

The above theorem assumes that  $N_1$  is divisible by  $S$ . If not, since  $N_1 \geq S$  (according to the discussion in Section 4.2.1) and  $S$  is typically small, we can always peel off some  $L_1$  iterations to make it divisible. In this way, we can assure at least the same performance as that of MS or xMS.

### 5.3 Relation to the Classical Modulo Scheduling of a Single Loop

If the loop nest is a single loop ( $n=1$ ), the sequential constraints are trivially satisfied. Other constraints are exactly the same as the those of the classical modulo scheduling. And the final schedule is  $f(o, (i_1)) = \sigma(o, i_1)$ . In this sense, classical MS is subsumed by SSP as a special case.

### 5.4 Relation to Hyperplane Scheduling

Next we establish the relation between our method and traditional hyperplane scheduling methods [8, 18]. We rewrite the mapping function for the final schedule in Equation (8) as follows.

$$f(o, \mathbf{I}) = \mathbf{I} \cdot \pi + offset(o, \mathbf{I}), \quad (16)$$

where  $\mathbf{I} = (i_1, i_2, \dots, i_n)$ , “ $\cdot$ ” is the inner product operator,

$$\pi = (T, (\prod_{2 < y \leq n+1} N_y) * S * T, \dots, N_{n+1} * S * T)^{transpose},$$

and

$$offset(o, \mathbf{I}) = \sigma(o, 0) + \underbrace{\left\lfloor \frac{i_1}{S} \right\rfloor * ((\prod_{2 \leq x \leq n+1} N_x) - 1) * S * T}_{(F)}. \quad (17)$$

The mapping function for the final schedule consists of two parts. The first part  $\mathbf{I}.\pi$  corresponds to hyperplane scheduling, which determines how to allocate the iteration points to slices. Unlike the traditional hyperplane scheduling that solves an integer programming problem to find out an optimal scheduling vector, here the scheduling vector  $\pi$  is predefined with  $S$  and  $T$  as parameters. The objective is thus not to find an optimal scheduling vector, but to find an optimal initiation interval.

This scheduling vector is unlikely to be found by the traditional hyperplane scheduling, because the vector expresses resource constraints through parameters  $S$  and  $T$ , while the traditional hyperplane scheduling does not consider resource constraints.

The second part,  $offset(o, \mathbf{I})$ , enforces dependences and resource constraints at the instruction level. In this offset, the first component is the 1-D schedule time,  $\sigma(o, 0)$ , which enforces those constraints within a slice, while the second component, (F), enforces resource constraints across slices. This offset is not a constant determined solely by the operation  $o$ ; it is a function of the first loop index  $i_1$ . Thus, the form of Equation (16) is similar to, but not a special case of, the known extensions of hyperplane scheduling [8, 23, 11]. The particular definition of this offset in Equation (17) is unlikely to be derived from the above methods.

## 5.5 Time Complexity

Our approach consists of loop selection, constructing a 1-D schedule, and computing the final schedule.

Loop selection is flexible and its complexity depends on the specific criteria. Let us consider the two criteria we presented in Section 4.2. Let  $UG$  be the number of uniformly generated sets, and  $u$  be the number of operations. In the worst case, for each loop level, we compute the lower bound of initiation interval that requires  $O(u^3)$  time with Floyd's All-Points Shortest Path algorithm [2], and estimate data reuse by Gauss-Seidal which requires  $O(UG * n^2)$  time. Therefore, the total time in this part is  $O(u^3 * n + UG * n^3)$ . In general,  $n$  is never greater than 6, and  $UG$  is typically small. Hence the dominant factor is still  $u$ . Thus the time complexity of the loop selection phase can be approximated as  $O(u^3)$ .

The construction of the 1-D schedule is traditional modulo scheduling applied to the simplified DDG, whose complexity is generally  $O(u^3)$  or  $O(u^4)$ , depending on the algorithm used [2]. The sequential constraints do not increase complexity. Computing the final schedule does not increase time, either, as it is simple parameter substitution.

To summarize, the overall time complexity of SSP is bounded by  $O(u^3)$  or  $O(u^4)$ , depending on the specific loop selection criteria and the modulo scheduling method used.

## 6 Extension to Imperfect Loop Nests

At instruction-level, it is common for a loop nest to be imperfect. Usually, a loop nest that is perfect in a high level representation becomes imperfect when lowered to instruction level. This is because operations for address calculation would be introduced between the loop levels.

Our study on scheduling perfect loop nests has set up a solid background, but cannot be applied directly to an imperfect loop nest: Not all operations appear with the same frequency now. An operation at an inner loop level runs more frequently than an operation at an outer loop level. Also, for efficiency, operations at different loop levels should be scheduled with different IIs, such that the instances of an operation at an inner loop level can be initiated at a faster rate, if possible.

In this section, we extend our scheduling approach to imperfect loop nests. First, we discuss how to schedule an imperfect loop nest with a single II. Subsequently, we discuss how the loop nest can be scheduled with multiple IIs to achieve higher execution efficiency.

## 6.1 Motivating Example

Fig. 8(a) shows an example imperfect loop nest. Compared with the example in Fig. 2, it differs only in the outer loop. Suppose we choose loop  $L_1$  for software pipelining. Let there be two function units, both being able to execute any operation, and each statement be considered as an operation with unit latency, except operation  $e$  with 2 cycles.

The effective DDG in Fig. 8(b) can be simplified to the 1-D DDG in Fig. 8(c), using the same concepts in Section 4. Based on this DDG and the resource constraints, we schedule all the operations as if they were in a single loop.

A 1-D schedule is shown in Fig. 8(d). The 1-D schedule is a *kernel nest* now: it has two kernels,  $K_1$  and  $K_2$ , corresponding to the two loops, and  $K_1$  encloses  $K_2$ . They have the same II. Stages 2, 3 and 4 contain the innermost loop operations. In general, we denote the total number of stages for the innermost loop as  $S_n$ . For this example,  $S_n = 3$ .

According to the 1-D schedule, ideally, all  $L_1$  iterations can be overlapped with the initiation interval of  $T = 3$  cycles, and each of them is sequential, as shown in Fig. 9(a). For uniformity of representation, we assume the operations before the innermost loop,  $a$ ,  $b$ , and  $c$ , have an index vector like  $(i_1, 0)$ , and the operation after the innermost loop,  $f$ , has an index vector like  $(i_1, N_2 - 1)$ .

In general, every  $S_n$  number of  $L_1$  iterations compose a group. After pushing down, we achieve the final schedule as illustrated in Fig. 9(b). To understand it, one may think the process as follows: at the beginning, an  $L_1$  iteration is issued every  $T$  cycles. After all the iterations in the first group have entered their innermost loop, they have filled the pipeline, and will hold the resources and continue running sequentially. All the other iterations stall in this period until the first group drains the pipeline and releases resources. At that time, they get resources and continue to issue. Such a process repeats until all iterations finish. Note that the draining of a previous group and the issuing of the next group are overlapped.

Such a way of execution leads to repeating patterns in the final schedule. Thus it can be rewritten into a compact form shown in Fig. 10. Like the perfect loop nest case, it is composed of a prolog, the repetition of an OLP and an ILES, and an epilog, except that the prolog is no longer a part of the first OLP. Each of them still consists of multiple copies of the kernel.

```

L1:for (i1=0; i1 < N1; i1++){
  a: W[i1]=1;
  b: U[i1][0]=W[i1];
  c: V[i1][0]=U[i1][0];
L2: for (i2=0; i2 < N2; i2++){
  d: U[i1 + 1][i2]=V[i1][i2]+U[i1][i2];
  e: V[i1][i2 + 1]=U[i1 + 1][i2];
  }
  f: W[i1]=V[i1][N2];
}

```

(a) An Example Imperfect Loop Nest

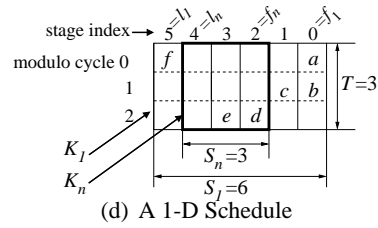
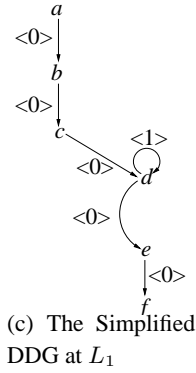
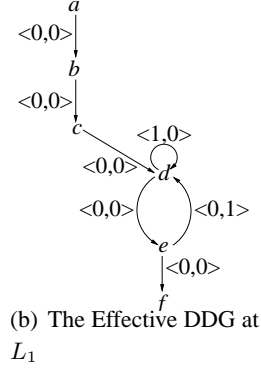


Figure 8: An Imperfect Loop Nest Example

## 6.2 Assumptions

We assume an imperfect loop nest model in Fig. 11, where each loop  $L_x$  has two sets of operations,  $OPSETA_x$  and  $OPSETB_x$ . For the innermost loop  $L_n$ ,  $OPSETA_n = OPSETB_n$ . For any loop  $L_x$ , its trip count  $N_x > 1$ .

If an operation  $o$  is in either  $OPSETA_x$  or  $OPSETB_x$ , it is said to be *at loop level*  $L_x$ , denoted as  $level(o) = x$ .

In general, an operation in  $OPSETA_x$  has an index of  $(i_1, i_2, \dots, i_x)$ . We can expand it to be an  $n$ -D vector  $(i_1, i_2, \dots, i_x, 0, \dots, 0)$  for convenience. Similarly, for an operation in  $OPSETB_x$ , we can expand its index to be an  $n$ -D vector  $(i_1, i_2, \dots, i_x, N_{x+1} - 1, \dots, N_n - 1)$ .

## 6.3 Solution with a Single Initiation Interval

Since any operation can be associated with an  $n$ -D index vector, a dependence distance vector is also an  $n$ -D vector; based on its value, the dependence can be classified as a zero, positive, or negative dependence. With this fact, the simplified DDG is defined in the same way as before.

Our approach remains to have the same three steps. Loop selection based on parallelism is the same as before. Loop selection based on data reuse can be the same as well if we estimate only the data reuse of the innermost loop, which is most frequently executed, and forget the operations at the outer loop

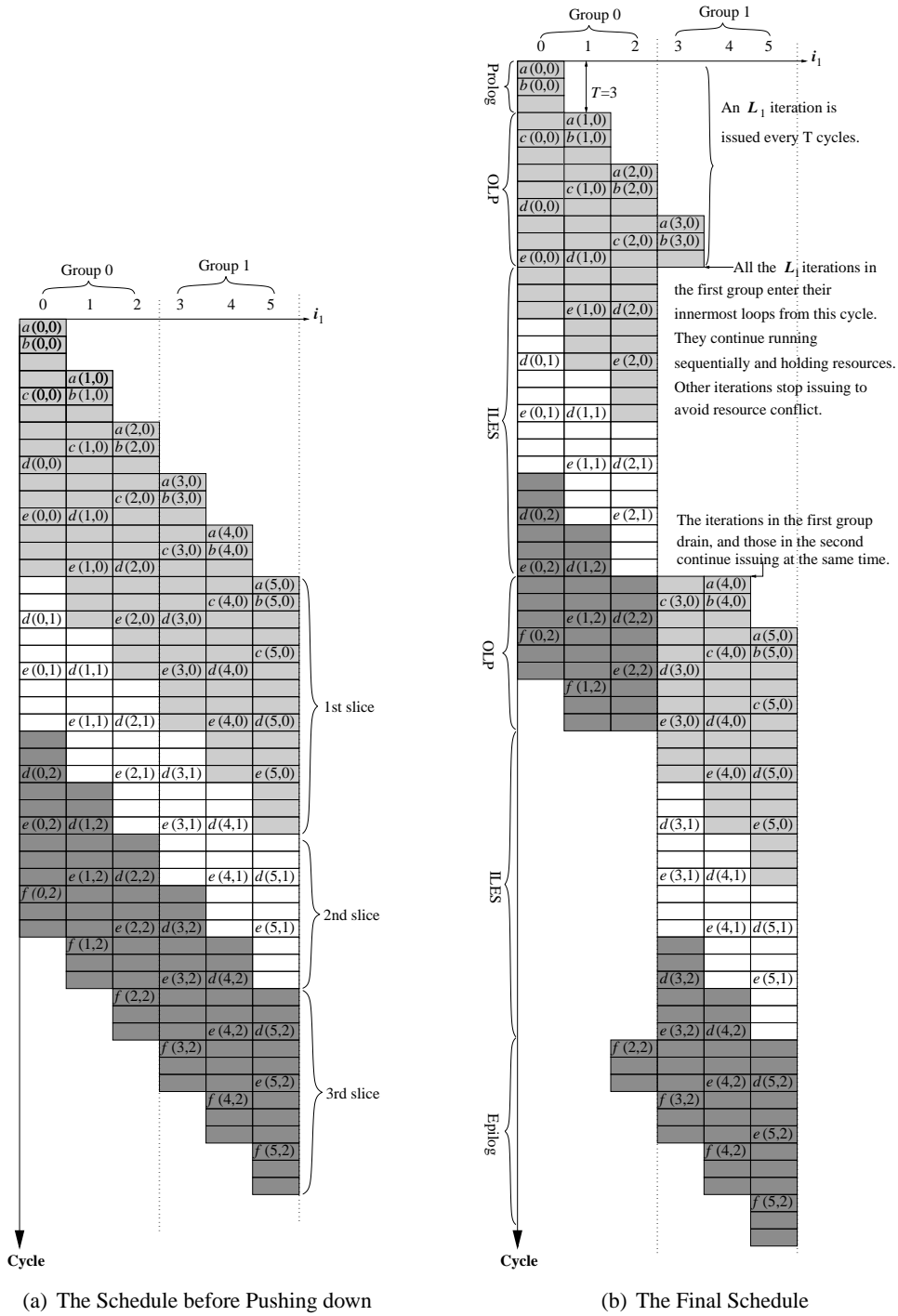


Figure 9: The Schedules before and after Pushing down (Assume  $N_1 = 6$  and  $N_2 = 3$ )

levels. Hence we focus on the next two steps, namely, 1-D schedule construction and final schedule computation. We assume the outermost loop  $L_1$  is chosen for software pipelining.

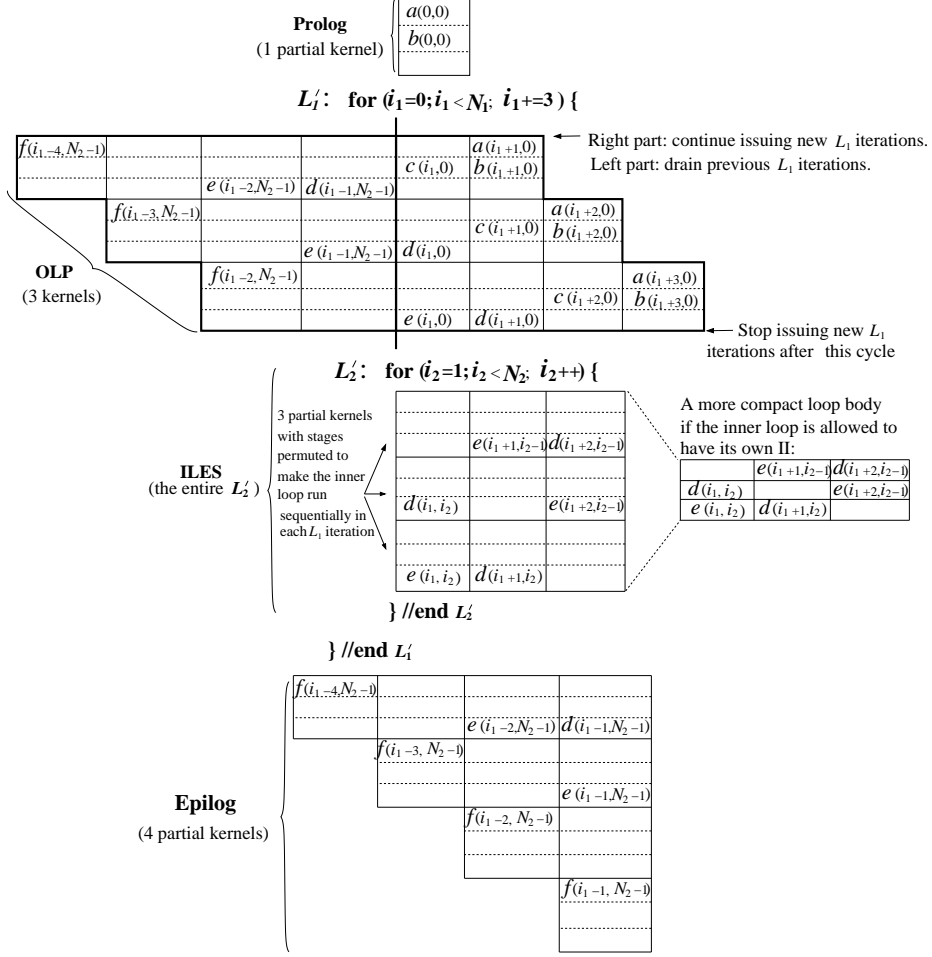


Figure 10: The Rewritten Loop Nest Representing the Final Schedule

### 6.3.1 1-D Schedule Construction

The 1-D schedule is now composed of  $n$  kernels, each corresponding to a loop (See Fig. 12). Kernel  $K_x$  is the kernel for loop  $L_x$ . Let  $f_x$  and  $l_x$  be the first and last stages of it, respectively. Then it has totally  $S_x = l_x - f_x + 1$  number of stages, including those of its inner loops. In general,  $S_n \leq S_{n-1} \leq \dots \leq S_2 \leq S_1$ . All kernels have the same initiation interval of  $T$  cycles.

Consistent with the nesting relationship of the loops,  $K_x$  contains  $K_{x+1}$ . Operations at an outer loop level are scheduled outside the stages of the inner loops.

As a convention, the 1-D schedule time is defined with the outermost loop kernel  $K_1$  as a reference. That is, for any operation  $o$ , if it is scheduled into modulo cycle  $q$  ( $0 \leq q < T$ ) in stage  $p$ , then its 1-D schedule time is  $\sigma(o, 0) = p * T + q$ . We also represent the stage as  $stage(o) = p$ .

The 1-D schedule needs to respect the following constraints:

1. Modulo property, dependence constraints, and resource constraints: they remain the same as those

```

L1: for (i1=0; i1 < N1; i1++) {
      OPSETA1
L2:   for (i2=0; i2 < N2; i2++) {
        OPSETA2
      ...
L_n:   ...   for (i_n=0; i_n < N_n; i_n++) {
            OPSETA_n
        } //end L_n
      ...
      OPSETB2
    } //end L2
  OPSETB1
} //end L1

```

Figure 11: The Imperfect Loop Nest Model

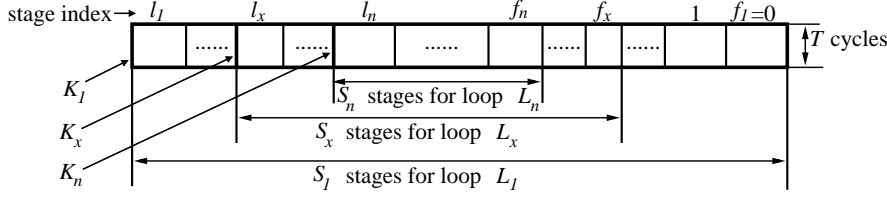


Figure 12: The 1-D Schedule with a Single II

in Section 4.3.

2. Sequential constraints: if  $n > 1$ , then for every positive dependence with  $o$  as the source operation,  $\delta$  being the dependence latency, and  $x = level(o)$ ,

$$\begin{cases} (l_x + 1) * T - \sigma(o, 0) \geq \delta & \text{if } x = n \text{ or } stage(o) > l_{x+1}, \\ f_{x+1} * T - \sigma(o, 0) \geq \delta & \text{otherwise.} \end{cases} \quad (18)$$

3. Kernel nesting constraints: (i)  $l_1 \geq l_2 \geq \dots \geq l_n \geq f_n \dots \geq f_2 \geq f_1$ , and (ii) for any operation  $o$ , if it is at loop level  $L_n$ , then  $stage(o) \in [f_n, l_n]$ . Otherwise, suppose it is at loop level  $L_x (x < n)$ , then  $stage(o) \in [f_x, f_{x+1}]$  if  $o$  is in  $OPSETA_x$ , or  $stage(o) \in (l_{x+1}, l_x]$  if  $o$  is in  $OPSETB_x$ .

The sequential constraints conservatively require operation  $o$  to complete before any possible use of it is issued, and thus the positive dependence from it must be respected.

The kernel nesting constraints express the nesting relationship of the kernels, and conservatively restrict the operations in  $OPSETA_x$  ( $OPSETB_x$ ) to be scheduled before (after) the inner loop kernel  $K_{x+1}$ .

### 6.3.2 Final Schedule Computation

For any operation  $o$  in an iteration point  $\mathbf{I}=(i_1, i_2, \dots, i_n)$ , the schedule time  $f(o, \mathbf{I})$  in Equation( 8) can be generalized as follows:

$$f(o, \mathbf{I}) = \sigma(o, i_1) + \sum_{2 \leq x \leq n} (i_x * ctime_x) + push(o, \mathbf{I}) * (ctime_1 - S_1 * T) \quad (19)$$

where  $ctime_x$  is the computation time of an  $L_x$  iteration in the ideal schedule where the outermost loop iterations are overlapped at the initiation interval of  $T$  cycles without delay, and  $push(o, \mathbf{I}) * (ctime_1 - S_1 * T)$  is the total delay that  $o(\mathbf{I})$  is pushed down to enforce resource constraints in the final schedule. In this delay,  $push(o, \mathbf{I})$  is the total number of ILESeS that appear before  $o(\mathbf{I})$  due to the pushing-down, and  $ctime_1 - S_1 * T$  is the length of an ILES<sup>5</sup>.

$Ctime_x$  can be recursively defined as

$$ctime_x = \begin{cases} (S_x - S_{x+1}) * T + N_{x+1} * ctime_{x+1} & \text{if } x < n, \\ S_n * T & \text{otherwise.} \end{cases} \quad (20)$$

The total ILESeS that appear before  $o(\mathbf{I})$  due to the pushing down is found to be as follows:

$$push(o, \mathbf{I}) = \begin{cases} \max(0, \lfloor \frac{i_1 + stage(o) - f_n + 1}{S_n} \rfloor) & \text{if } (i_2, \dots, i_n) = (0, \dots, 0) \\ & \text{and } stage(o) < f_n. \\ \min(\lfloor \frac{N_1 - 1}{S_n} \rfloor, \lfloor \frac{i_1 + stage(o) - l_n}{S_n} \rfloor) & \text{if } (i_2, \dots, i_n) = (N_2 - 1, \dots, N_n - 1) \\ & \text{and } stage(o) > l_n. \\ \lfloor \frac{i_1}{S_n} \rfloor & \text{otherwise.} \end{cases} \quad (21)$$

Let us briefly explain the definition. The operation instance  $o(\mathbf{I})$  is in group  $\lfloor \frac{i_1}{S_n} \rfloor$ . In general, the total number of ILESeS appearing before it due to the pushing-down is  $\lfloor \frac{i_1}{S_n} \rfloor$ . However, there are exceptions when  $o(\mathbf{I})$  is in the prolog, an OLP, or the epilog.

If  $o(\mathbf{I})$  is in the prolog, it is not pushed down at all. If it is in the right part of an OLP that fills new iterations, the total number of ILESeS equals  $\lfloor \frac{i_1 + stage(o) - f_n + 1}{S_n} \rfloor$ . Take  $a(5, 0)$  in Fig. 9(b) as an example. We have  $i_1 = 5$ ,  $stage(a) = 0$ ,  $f_n = 2$ , and  $S_n = 3$  according to the kernel in Fig. 8(d). Therefore, there is a total of  $\lfloor \frac{5+0-2+1}{3} \rfloor = 1$  ILES appearing before  $a(5, 0)$ . To summarize, we can simply express  $push(o, \mathbf{I})$  as  $\max(0, \lfloor \frac{i_1 + stage(o) - f_n + 1}{S_n} \rfloor)$ , which is the first case in Equation (21).

Second, if  $o(\mathbf{I})$  is in the epilog, the total number of ILESeS is  $\lfloor \frac{N_1 - 1}{S_n} \rfloor$ . If it is in the left side of an OLP that drains previous iterations, the number is  $\lfloor \frac{i_1 + stage(o) - l_n}{S_n} \rfloor$ . Take  $f(2, 2)$  in Fig. 9(b)

<sup>5</sup>It is easy to see from Fig. 9(b) or Fig. 10 that without the ILESeS, the final schedule is nothing but a modulo schedule, where each iteration has  $S_1 * T$  cycles. That implies that the length of an ILES equals  $ctime_1 - S_1 * T$  cycles.



as an example. We have  $i_1 = 2$ ,  $stage(f) = 5$ ,  $l_n = 4$ , and  $S_n = 3$  according to the kernel in Fig. 8(d). Thus the total number is  $\lfloor \frac{2+5-4}{3} \rfloor = 1$ . That is, it is delayed by one ILES, as can be seen from Fig. 9(b). Note that the ILES that delays it is the second ILES, not the first one, which is always before it and not due to the pushing-down, and thus not accounted for. In short, the total number is  $min(\lfloor \frac{N_1-1}{S_n} \rfloor, \lfloor \frac{i_1+stage(o)-l_n}{S_n} \rfloor)$ , which is the second case in Equation (21).

**Theorem 6.1.** *The final schedule defined in Equation (19) respects all the dependences in the effective DDG and the resource constraints.*

This theorem states the correctness of the final schedule. The proof is presented in the appendix.

### 6.3.3 Relation to the Scheduling of Perfect Loop Nests

When all the  $OPSETA_x$  and  $OPSETB_x$  ( $x < n$ ) are empty, the loop nest in Fig. 11 is perfect. Then  $l_1 = l_2 = \dots = l_n \geq f_n = \dots = f_2 = f_1$ ,  $S_n = S_{n-1} = \dots = S_2 = S_1$ , and all the operations are in the innermost loop. Consequently, the sequential constraints in Inequality (18) are equivalent to Inequality (7). And the kernel nesting constraints are trivially satisfied.

For the final schedule, since any stage is within  $[f_n, l_n]$ , we get  $push(o, \mathbf{I}) = \lfloor \frac{i_1}{S_n} \rfloor$ . Then the schedule time function defined in Equation (19) is equivalent to that in Equation (8).

In short, when the loop nest is perfect, both the 1-D schedule and the final schedule constructed by the method in this section are completely the same as those by the method in Section 4. In this sense, scheduling of an imperfect loop nest subsumes that of a perfect loop nest as a special case, as expected.

## 6.4 Solution with Multiple IIs

So far, we have assumed a single initiation interval for all the loop levels. To achieve better performance, however, it is desirable to have multiple IIs. Intuitively, the operations at an inner loop level run more frequently than those at an outer loop level, and therefore should run in a smaller II to shorten the execution time, whenever possible. This leads to the interesting topic of multi-II scheduling, which is also useful in practice.

Fig. 13 shows the general form of a 1-D schedule. Now each kernel  $K_x$  has its own initiation interval of  $T_x$  cycles. In general,  $T_1 \geq \dots \geq T_{n-1} \geq T_n$ . Although  $K_x$  takes only  $T_x$  cycles, the other cycles below and above it are empty, without any operation, as illustrated by the shaded places in the figure. We call them *null cycles*.

The 1-D schedule time is still defined with the outermost loop kernel  $K_1$  as a reference. That is, for any operation  $o$ , if it is scheduled into modulo cycle  $q$  ( $0 \leq q < T_1$ ) in stage  $p$ , then its 1-D schedule time is  $\sigma(o, 0) = p * T_1 + q$ .

The 1-D schedule needs to respect the following constraints, which are an extension of the constraints for the single-II case:

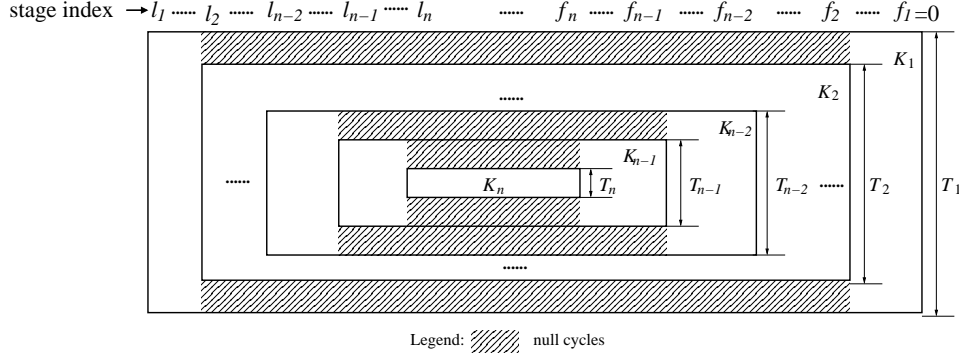


Figure 13: The General Form of the 1-D Schedule with Multiple IIs

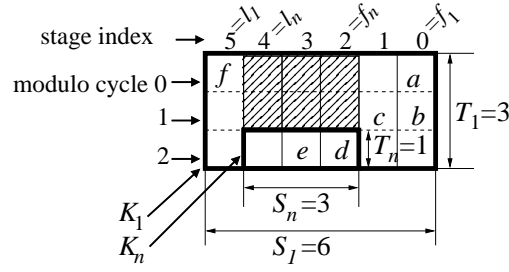


Figure 14: A 1-D Schedule with 2 IIs for the Loop Nest in Fig. 8(a). There are 2 null cycles above  $K_n$ .

1. Modulo property:

$$\sigma(o, i_1) + T_1 = \sigma(o, i_1 + 1) \quad (22)$$

2. Dependence constraints: for every dependence  $(o_1 \rightarrow o_2, \delta, \langle k \rangle)$  in the simplified DDG, let  $x = \min(\text{level}(o_1), \text{level}(o_2))$ ,  $y = \max(\text{level}(o_1), \text{level}(o_2))$ . Suppose  $\sigma(o_1, 0) = p_1 * T_1 + q_1$ , and  $\sigma(o_2, 0) = p_2 * T_1 + q_2$ , where  $0 \leq q_1, q_2 < T_1$ . Then

$$\begin{cases} k * T_y + (p_2 - p_1) * T_n + q_2 - q_1 \geq \delta & \text{if } \sigma(o_2, 0) \geq \sigma(o_1, 0), \\ k * T_y + (p_2 - p_1) * T_x + q_2 - q_1 \geq \delta & \text{otherwise.} \end{cases} \quad (23)$$

3. Resource constraints: At any modulo cycle in  $K_1$ , no hardware resource is allocated to more than one operation.
4. Sequential constraints: if  $n > 1$ , then for every positive dependence with  $o$  as the source operation,  $\delta$  being the dependence latency,

$$\begin{cases} (l_x + 1) * T_x - p * T_x - q \geq \delta & \text{if } x = n \text{ or } p > l_{x+1}, \\ f_{x+1} * T_x - p * T_x - q \geq \delta & \text{otherwise.} \end{cases} \quad (24)$$

where  $x = level(o)$ ,  $o$  is scheduled into modulo cycle  $q$  in stage  $p$ .

5. Kernel nesting constraints: (i)  $l_1 \geq l_2 \geq \dots \geq l_n \geq f_n \dots \geq f_2 \geq f_1$ , and  $T_1 \geq T_2 \geq \dots \geq T_n$   
(ii) for any operation  $o$ , if it is at loop level  $L_n$ , then  $stage(o) \in [f_n, l_n]$ . Otherwise, suppose it is at loop level  $L_x (x < n)$ , then  $stage(o) \in [f_x, f_{x+1}]$  if  $o$  is in  $OPSETA_x$ , or  $stage(o) \in (l_{x+1}, l_x]$  if  $o$  is in  $OPSETB_x$ .

The final schedule is hard to be described by a mapping function. It can be considered to be constructed in this way: first, use  $K_1$  to construct the final schedule as usual. That is, rewrite the loop nest into a parallel loop nest with  $K_1$  as the kernel. It is composed of  $n$  loops,  $L'_1, L'_2, \dots, L'_n$ . Each loop  $L'_x$  corresponds to the original loop  $L_x$ . Second, remove the null cycles. At the loop level of  $L'_x$ , only kernel  $K_x$  is involved. For example, in the final schedule in Fig. 10,  $L'_1$  is composed of  $K_1$ , and  $L'_2$  involves only  $K_2$  (The stages of this kernel permute, though). Therefore, the null cycles in the final schedule can be removed such that in  $L'_x$ , the kernel is “shrunk” from  $K_1$  into  $K_x$ , as illustrated by the inner loop in Fig. 10. In practice, the two steps can be combined: only the operations within the involved kernel is generated. The null cycles above it and below it are simply not produced. The dependences between the operations in this kernel and those outside it have been considered conservatively by the dependence and sequential constraints, such that without the null cycles, the dependences are still respected in the final schedule.

**EXAMPLE:** Fig. 14 shows an example kernel nest with two IIs for the loop nest in Fig. 8(a). With the outermost loop kernel  $K_1$  only, we have constructed a final schedule as shown in Fig. 10. Clearly, in  $L'_2$ , 2/3 of the total cycles are null cycles and unnecessarily wasted. After shrinking  $K_1$  to  $K_2$ , we reach a more compact schedule. See the annotation to the right of the figure. Now there are two initiation intervals: an outer loop iteration is issued at the II of 3 cycles, but after entering  $L'_2$ , an inner loop iteration is issued at the II of 1 cycle. The transition is natural without any special handling of the pipeline.

The correctness of the final schedule is shown in Appendix, which also contains an algorithm for constructing a 1-D schedule with multiple IIs, and a brief introduction to loop rewriting (code generation).

As expected, scheduling with multiple IIs subsumes scheduling with a single II as a special case. When all IIs are equal to a single value, say  $T$ , all the constraints for 1-D schedule construction become equivalent to those in Section 6.3. The two final schedules are of course the same, since their basic building blocks, the 1-D schedules, are the same and have no null cycles to remove.

## 7 Experiments

The SSP framework, including loop selection (with parallelism as the criterion), scheduling, register allocation, and code generation, was implemented in the ORC 2.1 compiler for the IA-64 architecture. The resulting code is run on an IA-64 Itanium workstation with a 733MHZ processor, 2GB main memory, and 16KB/96KB/2MB L1/L2/L3 caches, and the actual execution time is measured. Parallelism

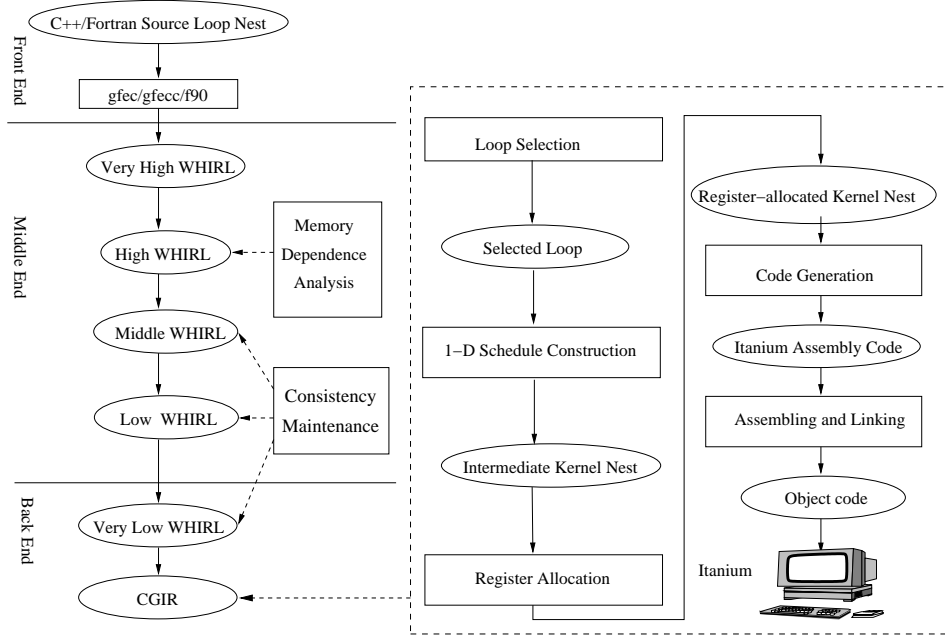


Figure 15: Compile Flow

serves as the first objective in loop selection, and ties are broken by data reuse, which is estimated manually by considering an abstract cache level with a line size of  $l$  according to Section 4.2.2.

Fig.15 shows the compile flow. The intermediate representation of the ORC compiler, termed as *WHIRL*, has 5 levels: Very High, High, Middle, Low, and Very Low levels, with increasing details and machine dependent information. The code generator translates the Very Low *WHIRL* to its own internal representation (CGIR) that matches the target machine instructions. Our implementation involves work from High *WHIRL* to CGIR.

At High level, it is relatively easy to get the multi-dimensional memory dependence information. This information is transmitted to the Very Low level throughout the process of Middle and Low levels, where many memory-related optimizations may happen. The information has to be consistently maintained during these optimizations.

At CGIR level, all instruction-level details have been exposed. The register dependences, and the memory dependences inherited from the High level, are combined together to build a complete DDG. Based on this DDG, a loop can be chosen either manually or automatically by estimating parallelism. Then for the selected loop, its  $n$ -D DDG is simplified to be 1-D. Based on the 1-D DDG and the underlying resource constraints, a 1-D schedule is constructed with multiple IIs. It is represented as a multiple-II kernel nest, or more exactly, a multiple-II *intermediate kernel nest*. It is intermediate in that variables in it have not been assigned registers yet. The register allocator works on this kernel nest and outputs a *register-allocated kernel nest*. From it, the code generator generates the IA-64 assembly code. The assembly code is then assembled, linked, and run on the Itanium machine.

The scheduling algorithm is left to the appendix. Register allocation and code generation are based on our previous work [Rong et al. 2005; Rong and Douillet et al. 2004], with minor extension to

accommodate our generic loop nest model in Fig. 11.

In this work, we report performance results for two loop kernels from scientific applications as well as loop nests extracted from SPEC2000 floating point benchmarks. The loop kernels that we consider are matrix multiplication (MM) and 2-D hydrodynamics (HD) modified from the Livermore Loops. We also apply three loop transformations, viz., loop interchange (6 different versions), loop tiling and unroll-and-jam, for MM and test each independently. The cache misses are measured for performance analysis using the IA64 performance monitoring tool, Pfmmon.

We extracted a number of loop nests from SPEC2000 floating point benchmarks. Each loop nest is wrapped as a function and called from the main routine with appropriate arguments. The function body, the loop nest, is compiled using our modified ORC compiler while the rest of the benchmark is compiled using gcc. This enables us to focus only on the implementation of SSP in the ORC compiler. The benchmarks are executed with the reference inputs of SPEC.

The loop level to be software pipelined by SSP can either be chosen by our compiler, or manually specified using a command line option. Let  $L_x$  be the loop level. We use  $SSP - L_x^*$  to represent the first case where  $L_x$  is chosen by our compiler, and  $SSP-L_x$  the second case where it is specified.

To verify the accuracy of our loop selection methods, SSP is applied to every feasible loop level of a loop nest. For each case, we compare the performance of the SSP-compiled loop nest (referred to as SSP), with that of modulo scheduling (referred to as MS), and that without software pipelining at all (referred to as Serial). Specially, for the tiled (unroll-and-jammed) MM, Serial refers to the original loop nest without tiling (unroll-and-jam) or software pipelining being applied, while MS and SSP refer to the software pipelined schedules after the loop nest is tiled (unroll-and-jammed). MS has been implemented in the original ORC distribution based on slack scheduling [14]. To test the effectiveness of SSP in the presence of other optimizations, the compiler optimization level is set to O3 (full optimization level) for all of Serial, MS, SSP.

Table 1 summarizes the average speedup for the loop nests tested. Speedup is defined as the execution time of a Serial loop nest divided by that of the optimized version (with MS, or with SSP applied to the selected loop level).

## 7.1 Performance of Kernel Loops

For MM, SSP always achieves the best speedup, with appropriate loop level being selected. Whether we apply loop interchange, tiling, or unroll-and-jam or no loop optimization at all, our method outperforms MS. Being able to work on a more profitable loop level, which is probably an outer loop level, allows the software-pipeliner to get around strong dependences or little data reuse opportunities of the innermost loop and to make use of the better properties of the other loops when possible. We discuss the performance in greater detail in the following subsections.

	MS	SSP	SSP over MS
MM- <i>ijk</i>	0.91	3.04	3.33
MM- <i>ikj</i>	0.94	4.5	4.79
MM- <i>jik</i>	0.98	3.5	3.57
MM- <i>jki</i>	0.88	4.43	5.02
MM- <i>kij</i>	0.97	2.54	2.62
MM- <i>kji</i>	0.96	2.51	2.62
HD	1.14	1.22	1.07
MM- <i>jki</i> with loop tiling	2.14	4.29	2.01
MM- <i>jki</i> with unroll-and-jam	6.19	10.17	1.64
168.wupwise loop 1	0.91	0.88	0.97
171.swim loop 1	0.83	0.83	1.0
171.swim loop 2	1.37	1.04	0.76
171.swim loop 3	1.0	0.97	0.97
173.applu loop 1	1.15	1.77	1.54
173.applu loop 2	0.97	1.94	1.99
173.applu loop 3	1.16	1.66	1.43
173.applu loop 4	1.11	2.23	2.01
173.applu loop 5	1.11	2.49	2.24
173.applu loop 6	1.14	2.22	1.95
173.applu loop 7	1.26	1.70	1.35
173.applu loop 8	0.81	1.06	1.31
301.apsi loop 1	0.99	1.86	1.87
301.apsi loop 2	1.03	4.07	3.97

Table 1: Average speedups

### 7.1.1 Matrix Multiply

We run SSP and MS on all the permutations of the matrix-multiply loop nest. The loop body is  $A[i][j] += B[i][k] * C[k][j]$ . The order of the loops in the nest are referred to as *ijk*, *ikj*, *jik*, *jki*, *kij*, and *kji*. Each loop order has different parallelism and data reuse potential. The performance results are depicted in Figure 16. We show the speedups achieved by MS and SSP for different matrix sizes.

For *ijk* and *jik*, the innermost loop is constrained by a recurrence cycle which limits the efficiency of MS. Consequently, applying SSP to other loop levels clearly achieves better performance. The upward tendency of the performance curves by software pipelining loop  $L_1$  suggests that with the increase in the matrix size, the advantage of SSP’s ability to retain data reuse becomes more important. For *ikj* and *jki*, the limited data reuse potential of one of the matrix operand of the innermost loop prevents MS to run efficiently when the size of the matrix increases. However, by software pipelining the outermost loop, such restriction is avoided. For *kij* and *kji*, all the loop levels show less data reuse potential or

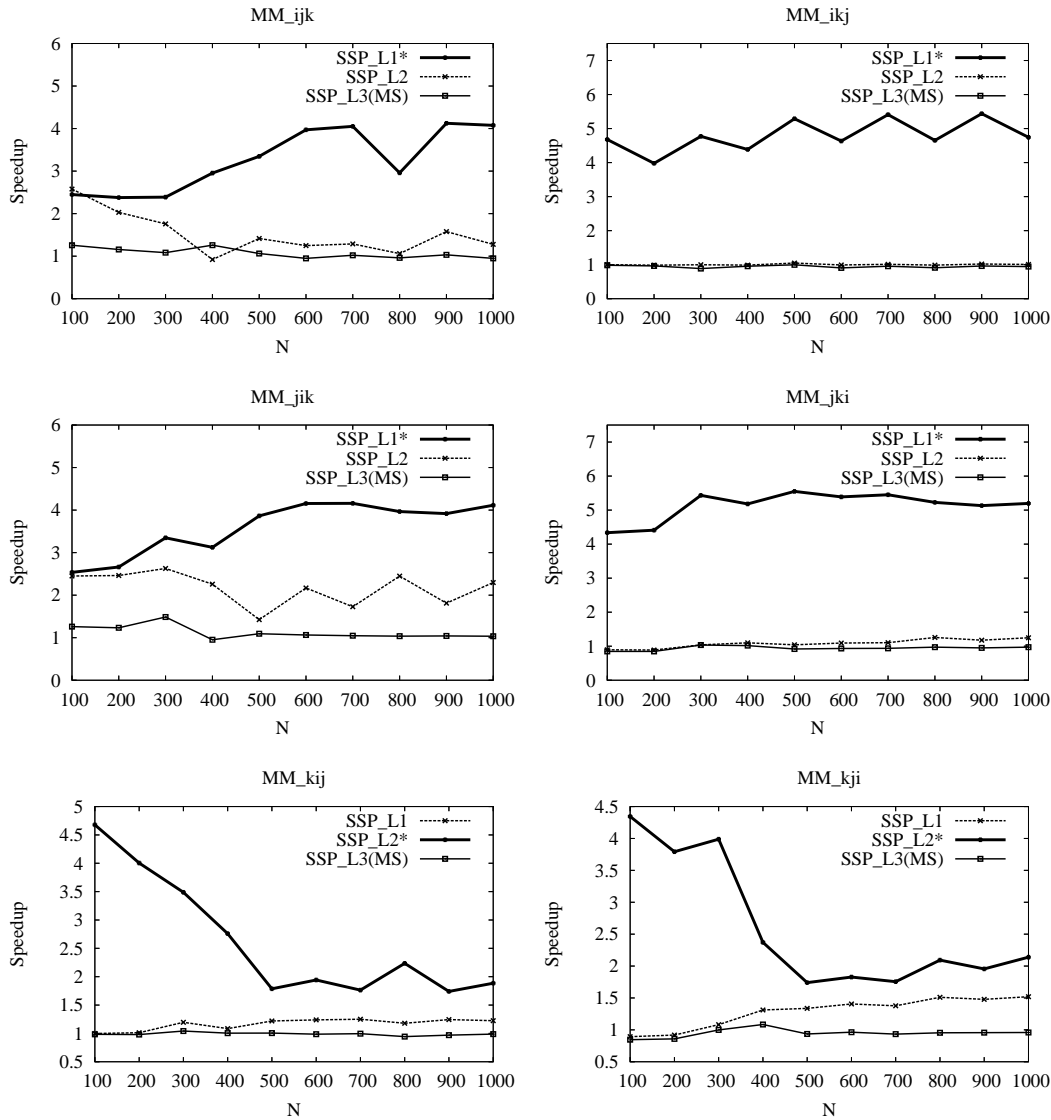


Figure 16: Performance of Matrix Multiply, where  $N * N$  is the Size of an Array

parallelism, limiting the speedup of all methods. However, software pipelining of the middle level still exhibits the best performance.

### 7.1.2 Tiled and Unroll-and-Jammed Matrix Multiply

Next we consider a classic tiled matrix multiply code [28] in which loop order is  $jki$ , and loops  $i$  and  $k$  are tiled <sup>6</sup>. The order of the loops is determined during tiling for best data locality. Hence only this order is considered in this experiment.

Figure 17 shows the performance for tiled matrix multiply with a constant tile size of 16 and the

<sup>6</sup>This is the equivalent row-major code. The column-major order is  $ikj$ , and loops  $k$  and  $j$  are tiled.

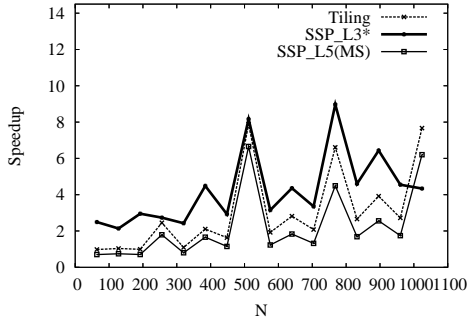


Figure 17: Tiled MM

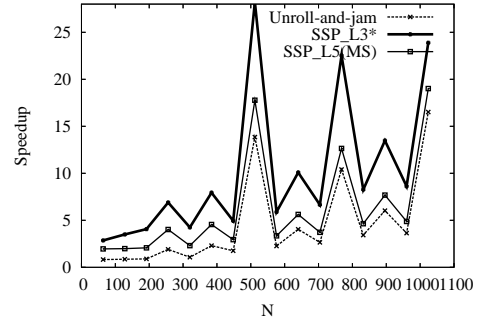


Figure 18: Unroll-And-Jammed MM

array sizes being multiples of it. After tiling, there is drastic improvement in speedup due to better data locality. The loop nest becomes 5-deep now. With software pipelining applied to the innermost loop, the performance is reduced by 38% on average. This is due to the overhead associated with the modulo schedule. Instead, when scheduling the third loop level, this overhead is amortized by the benefits from the longer execution time of the groups, and the natural overlapping of the draining and filling of adjacent groups.

Lastly, we consider an unroll-and-jammed version of MM (Figure 18). *Unroll-and-jam* [6], also known as *register tiling*, attempts to match the available parallelism in the application with the hardware resources. It is usually performed upon tiled code to further explore register level data reuse. Like in tiling, after unroll-and-jam, software pipelining of the middle loop level results in significant improvement.

The advantage of SSP scheduling shows clearly in these two experiments. Although the loops tested are perfect in high-level language, they become imperfect in assembly level. After tiling and unroll-and-jam, the depth of the whole loop nest becomes deeper (from 3 to 5), and the inner loops have small loop counts. And thus it becomes important to efficiently schedule the operations that are not in the innermost loop. It is also important to offset the overhead of initialization, finalization, and filling and draining the pipeline. Due to the small trip count of the innermost loop, such overheads have significant impact on the performance. By scheduling a middle loop level, software pipelining can offset these overheads effectively. The relatively longer execution time of a group outweighs the overhead. On the other hand, the operations at every loop level have been considered during 1-D schedule construction such that each loop level has the smallest possible initiation interval. In contrast, MS mainly cares about the efficiency of running the innermost loop operations and its software pipelined kernel includes only such operations.

### 7.1.3 Modified 2-D Explicit Hydrodynamics

The benchmark kernel considered is a 2-D explicit hydrodynamic code modified from Livermore loops. In this experiment, we varied the upper bounds  $kn$  and  $jn$ , respectively, of the outer and the inner loops.

Figure 19 shows the performance for the hydrodynamics benchmark when  $kn = jn$ . Since there is no recurrence in either loop level, data reuse will play a more important role in the performance. When



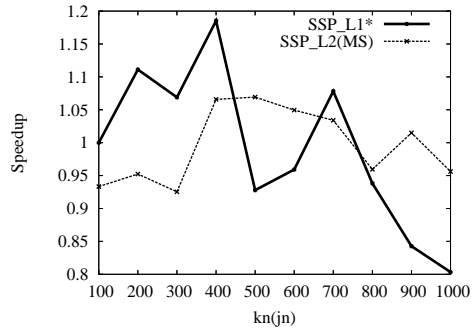


Figure 19: Hydrodynamics

the loop trip counts are smaller than 400, the outer loop is more beneficial. However, with an increase in the matrix size, the inner loop is better.

## 7.2 Performance of SPEC Loops

There are many loop nests in the SPEC2000 floating point benchmarks. However, many of them could not be software pipelined at an outer loop level due to either sibling inner loops, or function calls inside. Other loops have non-rectangular iteration spaces. We do not consider these loops.

In the loop nests extracted for experimentation, most of them have too short execution time and small loop counts (typically, less than 50 for a loop) to show meaningful performance improvement. However, they are perfectly fine for testing the correctness and effectiveness of our register allocation approach and heuristics, which we have reported in [Rong et al. 2005]. Here we report the performance of 14 loop nests from 168.wupwise, 171.swim, 173.applu, and 301.apsi, which have loop depths varying from 2 to 4 and have reasonable trip counts.

The results have been summarized in Table 1 for the appropriate chosen loop level. Here we study the results in more detail.

For 168.wupwise, there is a 3-deep critical loop nest. Scheduling the middle loop level results in a speedup of 0.97 over MS (i.e., a slowdown from 86.127 to 88.668 seconds). Software pipelining the outermost loop is possible, but register allocation method fails due to excessive integer register pressure.

For 171.swim, the table has shown the speedups for several loops. The least execution time within them is 32 seconds. Software pipelining the outer loop levels results in performance slowdown. However, in 173.applu, it improves performance by 30%-120%.

The two loop nests in 301.apsi have significant speedups when scheduling the outer loop levels: the first loop nest has good locality available in the outer loop, while the second loop nest has strong dependence cycles in the inner loop level.

The results suggest that any loop level, including the innermost one, may be the best or worst choice for software pipelining. Again, software pipelining should not be applied blindly to any loop level and loop selection is necessary, and the ability of SSP to exploit ILP from an arbitrary loop level is important.

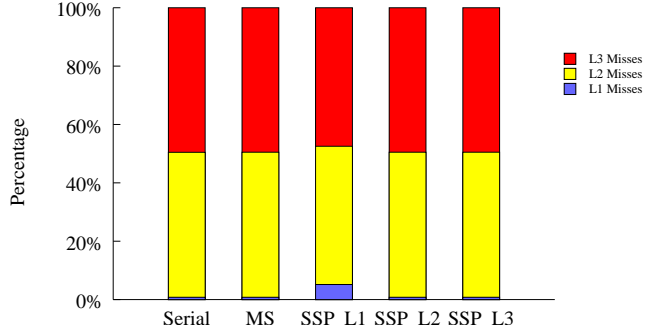


Figure 20: Percentages of the Three Level Cache Misses in MM- $ikj$  with Matrix Sizes 1024\*1024

### 7.3 Performance Analysis from Cache Misses

In this section, we investigate the cache effects of the different methods under comparison, with MM as an example. To correctly link the cache behavior to performance, we need to evaluate the relative weight of the cache misses at each cache level. Fig. 20 shows the relative weight of the cache misses at the 3 cache levels in the Itanium processor for  $ikj$  matrix multiply. Almost all cache misses happen at L2 and L3 caches for every method. This is the common trend in all the benchmarks tested.

This is a paradox as the first sight, since a cache miss at a low level must be caused by a cache miss at a high level, and thus L1 cache misses should be higher than L2. However, this does not hold for Itanium architecture. For this architecture: (1) There are separate L1 instruction cache and L1 data cache. And a floating point memory operation bypasses L1 data cache. Since our experiments are performed for scientific applications, and have only floating-point memory operations, L1 cache misses are mainly caused by instruction fetch. (2) The loop nests experimented are small in code size. Although SSP suffers from code expansion due to the lack of hardware support, the generated code is small enough to fit in the 16K bytes L1 instruction cache. Therefore, practically there are only a few instruction cache misses at the start of the loop nest.

From the above discussion, it is clear that we must focus on L2 and L3 cache misses in order to correctly explain the performance. For these experiments, we fix the matrix sizes as  $1024 \times 1024$ . The L2 and L3 cache misses (normalized with respect to the L2 and L3 cache misses of the serial code) are shown in Fig 21 and 22, where  $jk_{i+T}$  and  $jk_{i+UJ}$  refer to the tiled and unroll-and-jammed MM, respectively. They partly explain the performance improvement achieved by software pipelining a good loop level.

For all matrix multiply benchmarks, except the tiled and unroll-and-jammed versions, the L2 and L3 cache misses in SSP schedules at the outermost loop level are significantly lower than those for MS schedules. These numbers reflect that software pipelining exploits better data locality by selecting an appropriate loop level.

The increase in L2 and L3 cache misses in SSP schedules at the middle loop level  $L_3$  for the tiled and unroll-and-jammed benchmarks, at first, seems counter-intuitive as the SSP schedules have better speedups than the MS schedules for these benchmarks. However, in these two cases, MS affords cost in

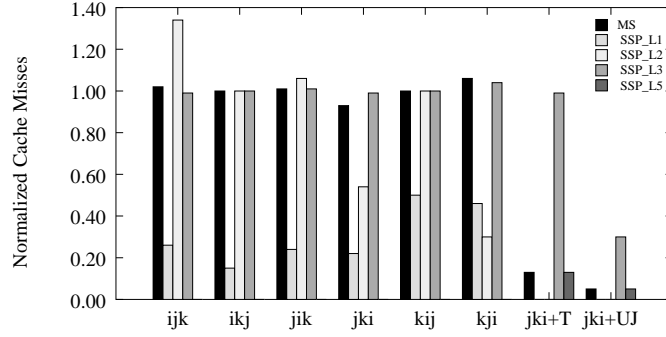


Figure 21: L2 Cache Misses

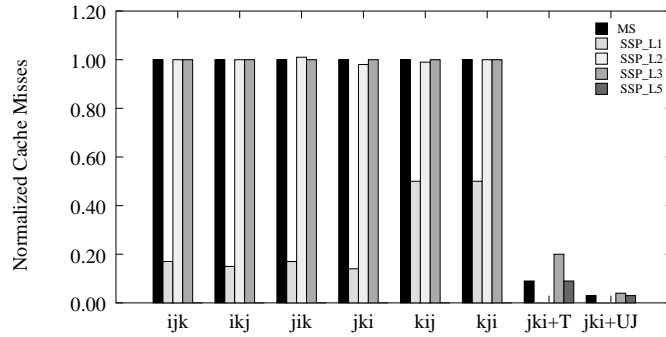


Figure 22: L3 Cache Misses

frequent pipeline filling and draining, and fails to schedule outer loop operations more effectively. The aggressive parallel execution of several iterations in SSP causes more cache pressure, leading to larger cache misses. The increase in cache misses affects the performance. Such effect is overcome partly by the increased parallelism in SSP schedules: when there are enough independent instructions that can be executed, the latencies due to cache misses can be masked. On the other hand, the Itanium architecture stalls only on uses. This also helps to overcome the effects of increased cache misses. The results also suggest that software pipelining tiled loops should consider the tile size to avoid negative cache misses, which we leave for future study.

## 8 Related Work

Most software pipelining algorithms [15, 2, 14, 25, 24] focus on the innermost loop, and do not consider cache effects. The most commonly used method, modulo scheduling, is a special case of our approach.

A common extension of modulo scheduling from single loops to loop nests, including *hierarchical reduction* [17], *Outer Loop Pipelining* [20], and *pipelining-dovetailing* [27], is to apply modulo scheduling hierarchically in order to exploit the parallelism between the draining and filling phases of adjacent outer loop iterations. In scheduling a loop, the DDG of its own is used.

In comparison, our method considers cache effects. The DDG is always the 1-D simplified DDG for the chosen loop, whatever loop level is currently under scheduling. The draining and filling phases

are naturally overlapped without any special treatment.

Loop Tiling [28] maximizes data locality, instead of parallelism. Loop unrolling duplicates the loop body of the innermost loop to increase instruction-level parallelism. Both methods are complementary to SSP.

One question is: *What is the difference between our method and the one that interchanges the selected loop with the innermost loop, and then software pipelines the new innermost loop with MS?* First, it may not always be possible to interchange the two loops. For example, if a dependence in a 3-deep loop nest has a distance vector of  $\langle 1, 1, -1 \rangle$  and our method selects the outermost loop, it is not legal to interchange this loop with the innermost loop. Second, even if they are interchangeable, the resulting schedules have different runtime performance due to different data reuse patterns. And for this interchanged loop nest, the choice for a good loop level might still be made by considering and comparing all the loop levels. Third, in some situations, interchange may be a bad choice, as we discussed in Section 3.1. Lastly, loop interchange can be beneficial to SSP as well.

Another question is: *What is the difference between our method and the one that tiles the selected loop, and then software pipelines the new innermost loop with MS?* In this case, the trip count of the new innermost loop is usually small as a result of tiling, and it is critical to hide the overhead of initialization, prolog, epilog, and finalization of the software pipelined innermost loop. Software pipelining an outer loop leads to less overhead, as discussed in Section 1 and confirmed in the experiments with tiled and unroll-and-jammed MM.

Unroll-and-jam [6] has been applied to improve the performance of software pipelined loops [5]. The outer loop is unrolled but it is still the innermost loop that is software-pipelined. The *RecMII* still strongly depends on the recurrences in the innermost loop, though reduced by the unroll factor. Unroll-and-squash first applies unroll-and-jam to a nested loop, and then reduce the code size of the jammed innermost loop by software pipelining and hardware support (rotating registers) [22]. SSP is different from unroll-and-squash in the following ways: (1) the unroll-and-squash method presented in [22] appears to be limited to 2-deep loop nest; (2) it does not overlap the epilog and prolog between successive outer loop iterations; and (3) it decides the unroll factor first, and then software pipelines the innermost loop.

In general, loop transformations, such as interchange, tiling, and unroll-and-jam, are orthogonal to our approach and can be applied independently. In Section 7, we have shown that our approach is beneficial with these loop transformations applied beforehand.

Hyperplane scheduling [18] is generally used in the context of large array-like hardware structures (such as systolic arrays and SIMD arrays), and does not consider resource constraints. There has been an interesting approach recently that enforces resource constraints to hyperplane scheduling by projecting the  $n$ -D iteration space to an  $(n - 1)$ -D virtual processor array, and then partitioning the virtual processors among the given set of physical processors [9]. This method targets parallel processor arrays, and does not consider low-level resources (like the function units within a single processor) or cache effects. A subsequent software pipelining phase may need to be applied to each physical processor in order to further exploit instruction-level parallelism from the iterations allocated to the same processor.

Other hyperplane-based methods [8, 23, 11] formulate the scheduling of loop nests as linear (often

integer linear) programming problems. Optimal solutions to integer programming have exponential time complexity in the worst case when using the Simplex algorithm or branch-and-bound methods [4]. Furthermore, they consider neither resource constraints nor cache effects.

Multi-dimensional retiming [21] translates a loop nest to be fully parallel without resource constraints.

Unimodular and non-unimodular transformations [4, 10] mainly care for coarse-grain parallelism or the communication cost between processors.

Fine-grain wavefront transformation [1] combines loop quantization and perfect pipelining to explore coarse and fine-grain parallelism simultaneously, based on outer loop unrolling and repetitive pattern recognition.

## 9 Conclusions

We have introduced the fundamental theory of software pipelining a loop nest at an arbitrary level that has a rectangular iteration space and has no sibling inner loops in it. This approach reduces the problem of  $n$ -dimensional software pipelining into a simpler problem of 1-dimensional software pipelining. This approach provides the freedom to search for and schedule the most profitable loop level, where profitability can be measured in terms of parallelism, data reuse potential, or any other criteria.

This approach subsumes the classical modulo scheduling as a special case. It also extends the traditional hyperplane scheduling to handle resource constraints. We have extended this approach to schedule imperfect loop nests. Multiple initiation intervals can be naturally achieved to issue operations at different loop levels in their fastest initiation rates.

We have demonstrated the correctness and efficiency of our method. Future work needs to study the interaction of this approach with other loop nest transformations like tiling, unroll-and-jam, and loop interchange more extensively, and extend the loop nest model to allow sibling inner loops at one level. Loop selection is another area to explore. This paper presented the principles in loop selection. In future, other data reuse models and other objectives like power consumption need to be investigated.

## ACKNOWLEDGMENTS

We are grateful to Prof. Bogong Su, Dr. Hongbo Yang, and the anonymous reviewers of the previous versions of this paper for their patient reviews and valuable advice.

## References

- [1] Alexander Aiken and Alexandru Nicolau. Fine-grain parallelization and the wavefront method. *Languages and Compilers for Parallel Computing*, MIT Press, Cambridge, MA, pages 1–16, 1990.

- [2] Vicki H Allan, Reese B Jones, Randall M Lee, and Stephen J Allan. Software pipelining. *ACM Computing Surveys*, 27(3):367–432, September 1995.
- [3] J. R. Allen, Ken Kennedy, Carrie Porterfield, and Joe Warren. Conversion of control dependence to data dependence. In *Conf. Record of 10th Annual ACM Symp. on Principles of Prog. Lang.*, pages 177–189, January 1983.
- [4] Utpal Banerjee. *Loop transformations for restructuring compilers: the foundations*. Kluwer Academic, 1993.
- [5] Steve Carr, Chen Ding, and Philip Sweany. Improving software pipelining with unroll-and-jam. In *Proc. 29th Annual Hawaii Intl. Conf. on System Sciences*, pages 183–192, 1996.
- [6] Steve Carr and Ken Kennedy. Improving the ratio of memory operations to floating-point operations in loops. *ACM Trans. on Prog. Lang. and Systems*, 16(6):1768–1810, Nov. 1994.
- [7] Steve Carr, Kathryn S. McKinley, and Chau-Wen Tseng. Compiler optimizations for improving data locality. In *Proceedings of the Sixth International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 252–262, San Jose, California, 1994.
- [8] Alain Darte and Yves Robert. Constructive methods for scheduling uniform loop nests. *IEEE Transactions on Parallel and Distributed Systems*, 5(8):814–822, August 1994.
- [9] Alain Darte, Robert Schreiber, B. Ramakrishna Rau, and Frédéric Vivien. Constructing and exploiting linear schedules with prescribed parallelism. *ACM Trans. Des. Autom. Electron. Syst.*, 7(1):159–172, 2002.
- [10] P. Feautrier. Automatic parallelization in the polytope model. *Lecture Notes in Computer Science*, 1132:79–103, 1996.
- [11] Guang R. Gao, Qi Ning, and Vincent Van Dongen. Software pipelining for nested loops. ACAPS Tech Memo 53, School of Computer Science, McGill Univ., Montréal, Québec, May 1993.
- [12] Somnath Ghosh, Margaret Martonosi, and Sharad Malik. Cache miss equations: a compiler framework for analyzing and tuning memory behavior. *ACM Transactions on Programming Languages and Systems*, 21(4):703–746, 1999.
- [13] R. Govindarajan, Erik R. Altman, and Guang R. Gao. A framework for resource-constrained rate-optimal software pipelining. *IEEE Trans. on Parallel and Distributed Systems*, 7(11):1133–1149, November 1996.
- [14] Richard A. Huff. Lifetime-sensitive modulo scheduling. In *Proc. of the ACM SIGPLAN '93 Conf. on Prog. Lang. Design and Implementation*, pages 258–267, Albuquerque, New Mexico, June 23–25, 1993. *SIGPLAN Notices*, 28(6), June 1993.
- [15] Intel. *Intel IA-64 Architecture Software Developer's Manual, Vol. 1: IA-64 Application Architecture*. Intel Corporation, Santa Clara, CA, USA, 2001.

- [16] Ken Kennedy and Kathryn S. McKinley. Optimizing for parallelism and data locality. In *ICS '92: Proceedings of the 6th international conference on Supercomputing*, pages 323–334, New York, NY, USA, 1992. ACM Press.
- [17] Monica Lam. Software pipelining: An effective scheduling technique for VLIW machines. In *Proc. Prog. Lang. Design and Impl.*, pages 318–328, June 22–24, 1988. *SIGPLAN Notices*, 23(7), July 1988.
- [18] L. Lamport. The parallel execution of DO loops. *Communications of the ACM*, 17(2):83–93, February 1974.
- [19] Soo-Mook Moon and Kemal Ebcioglu. Parallelizing nonnumerical code with selective scheduling and software pipelining. *ACM Transactions on Programming Languages and Systems*, 19(6):853–898, November 1997.
- [20] Kalyan Muthukumar and Gautam Doshi. Software pipelining of nested loops. *Lecture Notes in Computer Science*, 2027:165–181, 2001.
- [21] Nelson Luiz Passos and Edwin Hsing-Mean Sha. Achieving full parallelism using multidimensional retiming. *IEEE Trans. Parallel Distrib. Syst.*, 7(11):1150–1163, 1996.
- [22] D. Petkov, R. Harr, and S. Amarasinghe. Efficient pipelining of nested loops: unroll-and-squash. In *16th Intl. Parallel and Distributed Processing Symposium (IPDPS '02)*, Fort Lauderdale, FL, April 2002. IEEE.
- [23] J. Ramanujam. Optimal software pipelining of nested loops. In *Proc. of the 8th Intl. Parallel Processing Symp.*, pages 335–342, Cancún, Mexico, April 1994. IEEE.
- [24] B. R. Rau and J. A. Fisher. Instruction-level parallel processing: History, overview and perspective. *Journal of Supercomputing*, 7:9–50, May 1993.
- [25] B. Ramakrishna Rau. Iterative modulo scheduling: An algorithm for software pipelining loops. In *Proceedings of the 27th Annual International Symposium on Microarchitecture*, pages 63–74, San Jose, California, November 30–December 2, 1994.
- [26] Hongbo Rong, Alban Douillet, R. Govindarajan, and Guang R. Gao. Code generation for single-dimension software pipelining of multi-dimensional loops. In *CGO '04: Proceedings of the International Symposium on Code Generation and Optimization*, pages 175–186, Washington, DC, USA, 2004. IEEE Computer Society.
- [27] Jian Wang and Guang R. Gao. Pipelining-dovetailing: A transformation to enhance software pipelining for nested loops. volume 1060 of *Lecture Notes in Computer Science*, pages 1–17, Linköping, Sweden, April 1996.
- [28] Michael E. Wolf and Monica S. Lam. A data locality optimizing algorithm. In *Proc. of the ACM SIGPLAN '91 Conf. on Prog. Lang. Design and Implementation*, pages 30–44, Toronto, June 26–28, 1991. *SIGPLAN Notices*, 26(6), June 1991.

- [29] Michael E. Wolf, Dror E. Maydan, and Ding-Kai Chen. Combining loop transformations considering caches and scheduling. In *Proc. of the 29th Annual Intl. Symp. on Microarchitecture (MICRO 29)*, pages 274–286, Paris, December 2–4, 1996.



## A Appendix: Correctness, Scheduling Algorithm, and Loop Rewriting

In this section, we prove the correctness of the final schedule for an imperfect loop nest, first with a single II, then with multiple IIs. Then we describe the scheduling algorithm used in our prototype implementation for the experiments. Finally, we briefly illustrate how the loop nest is rewritten to realize the final schedule.

### A.1 Correctness of the Final Schedule with a Single II

First, let us establish the following facts:

**Lemma A.1.** *For any zero dependence  $(a \rightarrow b, \delta, \langle d_1, 0, \dots, 0 \rangle)$  in the effective DDG, if the dependence constraints are respected, then*

$$stage(b) + d_1 \geq stage(a). \quad (25)$$

The reason is as follows:

$$\begin{aligned} \sigma(b, d_1) - \sigma(a, 0) &\geq \delta \quad (\text{by Inequality (6)}) \\ &\geq 0 \quad (\delta \text{ is nonnegative as defined in Section 2}). \end{aligned}$$

Thus

$$\sigma(b, 0) + d_1 * T \geq \sigma(a, 0).$$

Dividing both sides by  $T$  and taking their floors, we reach the Inequality (25).

**Lemma A.2.** *For any zero dependence  $(a \rightarrow b, \delta, \langle d_1, 0, \dots, 0 \rangle)$  in the effective DDG, if the dependence constraints are respected, then*

$$stage(a) > l_n \iff stage(b) > l_n$$

According to the kernel nesting constraints in Section 6.3.1,  $stage(a) > l_n \iff$  operation  $a$  is in some  $OPSETB_x$ , and therefore its index  $\mathbf{I}$  is in the form of  $(i_1, i_2, \dots, i_x, N_{x+1} - 1, \dots, N_n - 1)$ . Then  $\mathbf{I}'$  is in the form of  $(i_1 + d_1, i_2, \dots, i_x, N_{x+1} - 1, \dots, N_n - 1)$ . Only operations in  $OPSETB_x$  could have this form of index vector. Therefore  $stage(b) > l_n$  as well.

**Lemma A.3.** *For any operation  $o$ , and two index vectors  $\mathbf{I}$  and  $\mathbf{I}'$ , if  $\mathbf{I}' - \mathbf{I} = (d_1, 0, \dots, 0)$ , where  $d_1 \geq 0$ , then*

$$push(o, \mathbf{I}) \leq push(o, \mathbf{I}')$$

The proof is trivial based on the definition of function  $push$ . From this fact, we can easily see that  $f(o, \mathbf{I}) \leq f(o, \mathbf{I}')$ , and the equality holds only when  $d_1$  is 0. This confirms our intuition that the instances of the same operation from successive outermost loop iterations, and within the same slice before the pushing down, are still issued in the same order after that (The only change is the distances between them, not the order).

**Lemma A.4.** For any operation  $o$ , and three index vectors:  $\mathbf{I}$ ,  $\mathbf{I}'$ , and  $\mathbf{I}''$ , if  $\mathbf{I} = (i_1, 0, \dots, 0)$ ,  $\mathbf{I}' = (i_1, N_2 - 1, \dots, N_n - 1)$ , and  $\mathbf{I}'' = (i_1, i_2, \dots, i_n)$  is not equal to  $\mathbf{I}$  nor  $\mathbf{I}'$ , where  $i_1$  is any value for the first index, then

$$\text{push}(o, \mathbf{I}) \leq \left\lfloor \frac{i_1}{S_n} \right\rfloor = \text{push}(o, \mathbf{I}') \leq \text{push}(o, \mathbf{I}'')$$

*Proof.* Under the given conditions,

$$\begin{aligned} \text{push}(o, \mathbf{I}) &= \begin{cases} \max(0, \left\lfloor \frac{i_1 + \text{stage}(o) - f_n + 1}{S_n} \right\rfloor) & \text{if } \text{stage}(o) < f_n \\ \left\lfloor \frac{i_1}{S_n} \right\rfloor & \text{otherwise.} \end{cases} \\ \text{push}(o, \mathbf{I}') &= \left\lfloor \frac{i_1}{S_n} \right\rfloor \\ \text{push}(o, \mathbf{I}'') &= \begin{cases} \min\left(\left\lfloor \frac{N_1 - 1}{S_n} \right\rfloor, \left\lfloor \frac{i_1 + \text{stage}(o) - l_n}{S_n} \right\rfloor\right) & \text{if } \text{stage}(o) > l_n \\ \left\lfloor \frac{i_1}{S_n} \right\rfloor & \text{otherwise.} \end{cases} \end{aligned} \quad (26)$$

In the first case of  $\text{push}(o, \mathbf{I})$ ,  $i_1 + \text{stage}(o) - f_n + 1 \leq i_1$ . Therefore,  $\text{push}(o, \mathbf{I})$  can not be more than  $\left\lfloor \frac{i_1}{S_n} \right\rfloor$ .

Similarly, in the first case of  $\text{push}(o, \mathbf{I}'')$ ,  $N_1 - 1 \geq i_1$  and  $i_1 + \text{stage}(o) - l_n > i_1$ . Therefore,  $\text{push}(o, \mathbf{I}'')$  cannot be less than  $\left\lfloor \frac{i_1}{S_n} \right\rfloor$ .  $\square$

The above lemma confirms the intuition that within the same outermost loop iteration, with the increase of the index vectors, an operation can only be pushed further from its original position.

**Corollary A.5.** Given any two operations  $o$  and  $o'$ , and their index vectors  $\mathbf{I}$  and  $\mathbf{I}'$ , if their distance  $\mathbf{I}' - \mathbf{I} = \langle 0, d_2, \dots, d_n \rangle$  is positive, then

$$\text{push}(o, \mathbf{I}) \leq \text{push}(o', \mathbf{I}'). \quad (27)$$

*Proof.* Let  $\mathbf{I} = (i_1, i_2, \dots, i_n)$ . Then  $\mathbf{I}' = (i_1, i_2 + d_2, \dots, i_n + d_n)$ .

Because the difference between the two index vectors is positive,  $\mathbf{I}$  cannot be in this form:  $(i_1, N_2 - 1, \dots, N_n - 1)$ . Therefore, according to the proof of Lemma A.4, we know

$$\text{push}(o, \mathbf{I}) \leq \left\lfloor \frac{i_1}{S_n} \right\rfloor.$$

For the same reason,  $\mathbf{I}'$  is impossible to be in this form:  $(i_1, 0, \dots, 0)$ , and by the proof of Lemma A.4, there is

$$\left\lfloor \frac{i_1}{S_n} \right\rfloor \leq \text{push}(o', \mathbf{I}').$$

$\square$

Note that  $o(\mathbf{I})$  and  $o'(\mathbf{I}')$  are within the same outermost loop iteration, and the two operations are arbitrary. This corollary says that the final schedule is consistent with our principle that an outermost loop iteration is executed sequentially.

**Lemma A.6.** *Given a positive dependence with the distance vector of  $\langle d_1, \dots, d_n \rangle$ , and  $d_y$  is the first non-zero (actually positive) element in the sub-vector  $\langle d_2, \dots, d_n \rangle$ , then*

$$\sum_{2 \leq x \leq n} (d_x * ctime_x) \geq S_y * T$$

*Proof.* The sum is minimum when  $d_y$  takes the smallest positive number (1), and any other element after it takes the smallest legal distance. Since an index  $i_z$  for loop  $L_z$  is between  $[0, N_z - 1]$ , the smallest legal distance corresponding to this loop is  $-N_z + 1$ . So the minimum sum is

$$\begin{aligned} ctime_y + \sum_{y < x \leq n} ((-N_x + 1) * ctime_x) &= \underbrace{ctime_y - N_{y+1} * ctime_{y+1}}_{(A)} + \\ &\quad \underbrace{ctime_{y+1} - N_{y+2} * ctime_{y+2}}_{(B)} + \\ &\quad \dots \\ &\quad \underbrace{ctime_{n-1} - N_n * ctime_n}_{(C)} + \\ &\quad ctime_n \end{aligned}$$

From the definition of the computation time of a loop,

$$ctime_y = (S_y - S_{y+1}) * T + N_{y+1} * ctime_{y+1}.$$

So

$$(A) = (S_y - S_{y+1}) * T.$$

Similarly,

$$(B) = (S_{y+1} - S_{y+2}) * T,$$

$\dots,$

$$(C) = (S_{n-1} - S_n) * T.$$

It is clear that the minimum sum is  $S_y * T - S_n * T + ctime_n = S_y * T$ .  $\square$

Now we are ready to prove the correctness of the final schedule for an imperfect loop nest with a single II. The correctness is established by Theorem 6.1 in Section 6.3.2.

*Proof.* Given a dependence  $(a \rightarrow b, \delta, \langle d_1, d_2, \dots, d_n \rangle)$ ,  $\mathbf{I} = (i_1, i_2, \dots, i_n)$ , and  $\mathbf{I}' = (i_1 + d_1, i_2 + d_2, \dots, i_n + d_n)$ , there will be

$$\begin{aligned}
f(b, \mathbf{I}') - f(a, \mathbf{I}) &= \underbrace{\sigma(b, i_1 + d_1) - \sigma(a, i_1)}_{(A)} + \\
&\quad \underbrace{\sum_{2 \leq x \leq n} (d_x * ctime_x)}_{(B)} + \\
&\quad \underbrace{(push(b, \mathbf{I}') - push(a, \mathbf{I})) * (ctime_1 - S_1 * T)}_{(C)}. \tag{28}
\end{aligned}$$

1. If the dependence is a zero dependence, then

$$\begin{aligned}
(A) &\geq \delta \quad (\text{By the dependence constraints in Inequality (6)}), \\
(B) &= 0.
\end{aligned}$$

In order for the dependence to be respected in the final schedule, i.e.,  $(A) + (B) + (C) \geq \delta$ , we need to prove

$$(C) \geq 0.$$

According to our assumption in Section 6.2,  $(N_2 - 1, \dots, N_n - 1) \neq (0, \dots, 0)$ . Then it is easy to know that  $ctime_1 > S_1 * T$ . Therefore, for  $(C) \geq 0$ , we need only to prove

$$push(b, \mathbf{I}') - push(a, \mathbf{I}) \geq 0.$$

Below we prove it by enumerating all possibilities.

Case 1:  $(i_2, \dots, i_n)$  is not equal to  $(0, \dots, 0)$ , nor  $(N_2 - 1, \dots, N_n - 1)$ . Then

$$push(b, \mathbf{I}') - push(a, \mathbf{I}) = \left\lfloor \frac{i_1 + d_1}{S_n} \right\rfloor - \left\lfloor \frac{i_1}{S_n} \right\rfloor \geq 0.$$

Case2:  $(i_2, \dots, i_n) = (0, \dots, 0)$ . Then

$$\begin{aligned}
push(b, \mathbf{I}') &= \begin{cases} \max(0, \left\lfloor \frac{i_1 + d_1 + stage(b) - f_n + 1}{S_n} \right\rfloor) & \text{if } stage(b) < f_n \\ \left\lfloor \frac{i_1 + d_1}{S_n} \right\rfloor & \text{otherwise.} \end{cases} \\
push(a, \mathbf{I}) &= \begin{cases} \max(0, \left\lfloor \frac{i_1 + stage(a) - f_n + 1}{S_n} \right\rfloor) & \text{if } stage(a) < f_n \\ \left\lfloor \frac{i_1}{S_n} \right\rfloor & \text{otherwise.} \end{cases}
\end{aligned}$$

$Push(a, \mathbf{I}) \leq push(b, \mathbf{I}')$  because when

$$stage(a) < f_n, \tag{29}$$

there is

$$\left\lfloor \frac{i_1 + stage(a) - f_n + 1}{S_n} \right\rfloor \leq \begin{cases} \left\lfloor \frac{i_1 + d_1 + stage(b) - f_n + 1}{S_n} \right\rfloor & \text{(by Lemma A.1),} \\ \left\lfloor \frac{i_1 + d_1}{S_n} \right\rfloor & \text{(by Inequality (29)).} \end{cases}$$

While when

$$stage(a) \geq f_n, \quad (30)$$

there is

$$\left\lfloor \frac{i_1}{S_n} \right\rfloor \leq \begin{cases} \left\lfloor \frac{i_1 + stage(a) - f_n + 1}{S_n} \right\rfloor \leq \left\lfloor \frac{i_1 + d_1 + stage(b) - f_n + 1}{S_n} \right\rfloor & \text{(by Inequality (30) and Lemma A.1)} \\ \left\lfloor \frac{i_1 + d_1}{S_n} \right\rfloor. \end{cases}$$

Case 3:  $(i_2, \dots, i_n) = (N_2 - 1, \dots, N_n - 1)$ . Let us examine the two cases of operation  $a$ .

When

$$stage(a) > l_n,$$

we have

$$stage(b) > l_n \text{ (by Lemma A.2),}$$

and thus

$$\begin{aligned} push(b, \mathbf{I}') &= \min\left(\left\lfloor \frac{N_1 - 1}{S_n} \right\rfloor, \left\lfloor \frac{i_1 + d_1 + stage(b) - l_n}{S_n} \right\rfloor\right) \\ push(a, \mathbf{I}) &= \min\left(\left\lfloor \frac{N_1 - 1}{S_n} \right\rfloor, \left\lfloor \frac{i_1 + stage(a) - l_n}{S_n} \right\rfloor\right). \end{aligned}$$

$Push(a, \mathbf{I}) \leq push(b, \mathbf{I}')$  as

$$\left\lfloor \frac{i_1 + stage(a) - l_n}{S_n} \right\rfloor \leq \left\lfloor \frac{i_1 + d_1 + stage(b) - l_n}{S_n} \right\rfloor \text{ (by Lemma A.1).}$$

When

$$stage(a) \leq l_n,$$

we have

$$stage(b) \leq l_n \text{ (by Lemma A.2),}$$

and thus

$$\begin{aligned} push(b, \mathbf{I}') &= \left\lfloor \frac{i_1 + d_1}{S_n} \right\rfloor, \\ push(a, \mathbf{I}) &= \left\lfloor \frac{i_1}{S_n} \right\rfloor. \end{aligned}$$

Obviously,  $push(a, \mathbf{I}) \leq push(b, \mathbf{I}')$ .

2. If the dependence is a positive dependence, let  $x = level(a)$ .

Case 1:  $x = n$  or  $stage(a) > l_{x+1}$ . That is, operation  $a$  is in  $L_n$  or in a  $OPSETB_x$  of an outer loop, according to the kernel nesting constraints.

Imagine there is an operation  $c$  in the last stage of the same loop, stage  $l_x$ , and at the last modulo cycle of it,  $T - 1$ . Its 1-D schedule time is then

$$\sigma(c, 0) = (l_x + 1) * T - 1 \quad (31)$$

Whether such an operation really exist does not matter: we only want to take it as a reference. In order to prove  $f(b, \mathbf{I}') - f(a, \mathbf{I}) \geq \delta$ , we only need to prove the following facts:

$$f(c, \mathbf{I}) - f(a, \mathbf{I}) \geq \delta - 1 \quad (32)$$

$$f(b, \mathbf{I}') - f(c, \mathbf{I}) \geq 1 \quad (33)$$

We know

$$\begin{aligned} f(c, \mathbf{I}) - f(a, \mathbf{I}) &= \underbrace{\sigma(c, i_1) - \sigma(a, i_1)}_{(D)} + \\ &\quad \underbrace{\sum_{2 \leq x \leq n} (0 * ctime_x)}_{(E)} + \\ &\quad \underbrace{(push(c, \mathbf{I}) - push(a, \mathbf{I})) * (ctime_1 - S_1 * T)}_{(F)} \end{aligned}$$

It is easy to see that

$$(D) = (l_x + 1) * T - 1 - \sigma(a, 0) \geq \delta - 1 \text{ (by Equation ( 31) and Inequality ( 18))},$$

$$(E) = 0,$$

$$(F) \geq 0.$$

Thus Inequality ( 32) holds.

Now we prove Inequality ( 33) is true.

$$\begin{aligned} f(b, \mathbf{I}') - f(c, \mathbf{I}) &= \underbrace{\sigma(b, d_1) - \sigma(c, 0)}_{(G)} + \\ &\quad \underbrace{\sum_{2 \leq x \leq n} (d_x * ctime_x)}_{(H)} + \\ &\quad \underbrace{(push(b, \mathbf{I}') - push(c, \mathbf{I})) * (ctime_1 - S_1 * T)}_{(I)} \end{aligned}$$

We will prove that  $(G) + (H) \geq 1$  and  $(I) \geq 0$ .

We know

$$(G) \geq \sigma(b, 0) - \sigma(c, 0) \geq \sigma(b, 0) - (l_x + 1) * T + 1 \text{ (by Equation ( 31))}, \quad (34)$$

$$(H) \geq S_y * T \text{ (by Lemma A.6)}, \quad (35)$$

where  $y$  is such that  $d_y$  is the first non-zero (actually positive) element in the sub-vector  $\langle d_2, \dots, d_n \rangle$ . Remember that operation  $c$  (and  $a$ ) is in  $OPSETB_x$  and has an index vector of  $\mathbf{I} = \langle i_1, \dots, i_x, N_{x+1} - 1, \dots, N_n - 1 \rangle$ . It is not possible to make the elements after  $i_x$  bigger. It is only possible to make one or more of the elements between  $i_1$  and  $i_x$  bigger. This implies that the first non-zero increment,  $d_y$ , must be applied to one of the elements between  $i_1$  and  $i_x$ . In other words,  $y \leq x$ , and thus

$$l_y \geq l_x \quad (36)$$

Operation  $b$  must be at loop level  $L_y$  or any inner loop level of it. If otherwise, it is outside of this loop, its index corresponding to this loop would be either 0 or  $N_y - 1$ , depending on whether it is in an  $OPSETA$  or  $OPSETB$ , and cannot be changed, let alone a non-zero increment  $d_y$ . Therefore,

$$\sigma(b, 0) \geq f_y * T. \quad (37)$$

Now let us add up the Inequalities( 34) and ( 35):

$$\begin{aligned} (G) + (H) &\geq f_y * T - (l_x + 1) * T + 1 + S_y * T \text{(by Inequalities ( 37))} \\ &= (l_y - l_x) * T + 1 \\ &\geq 1 \text{(by Inequality ( 36))}. \end{aligned} \quad (38)$$

Now we need only to prove that  $(I) \geq 0$ . Let an index vector  $\mathbf{I}'' = \langle i_1 + d_1, i_2, \dots, i_n \rangle$ . That is, it is the same as  $\mathbf{I}$  except for the first element. We decompose  $(I)$  as follows:

$$(I) = \underbrace{push(b, \mathbf{I}') - push(c, \mathbf{I}'')}_{(J)} + \underbrace{push(c, \mathbf{I}'') - push(c, \mathbf{I})}_{(K)}. \quad (39)$$

$(K) \geq 0$  according to Lemma A.3. There must be  $(J) \geq 0$  as well, because  $\mathbf{I}' - \mathbf{I}'' = \langle 0, d_2, \dots, d_n \rangle$  is a positive vector. According to corollary A.5,  $(J)$  must be nonnegative.

Case 2: Otherwise, similarly, imagine there is an operation  $c$  at the last modulo cycle before the inner loop kernel. That is, the 1-D schedule time of it is

$$\sigma(c, 0) = f_{x+1} * T - 1 \quad (40)$$

We still prove the facts stated in Inequalities ( 32) and ( 33). The proof is the same as that for Case 1, except the following details:

$$\begin{aligned} (G) &\geq \sigma(b, 0) - f_{x+1} * T + 1 \text{ (by Equation ( 40))}, \\ (G) + (H) &\geq f_y * T - f_{x+1} * T + 1 + S_y * T \\ &= (l_y - f_{x+1} + 1) * T + 1 \\ &> 1 \end{aligned}$$

### 3. Resource constraints.

Lastly, for any two distinct operation instances,  $a(\mathbf{I})$  and  $b(\mathbf{I}')$ , scheduled at the same cycle, we have

$$f(b, \mathbf{I}') - f(a, \mathbf{I}) = 0.$$

Borrow Equation (28), we have  $(A) + (B) + (C) = 0$ . Since  $(A) = \sigma(b, i_1) + d_1 * T - \sigma(a, i_1)$ , and  $(B) + (C)$  is a multiple of  $T$ , we have  $\sigma(b, i_1) - \sigma(a, i_1) = -d_1 * T - (B) - (C)$ , which is also a multiple of  $T$ . That is, operations  $a$  and  $b$  must be from the same modulo cycle in the 1-D schedule, which has no resource conflicts since the 1-D schedule respects the resource constraints.

□

## A.2 Correctness of the Final Schedule with Multiple IIs

First, suppose all IIs are equal to a single value, say,  $T$ . Then all the scheduling constraints, and thus the final schedule, are equivalent to those with a single II, as discussed in Section 6.4. We have proved the correctness of such a final schedule in the above section.

Second, from this final schedule, let us remove all unnecessary “null cycles”, and reach a multiple-II final schedule. This final schedule can be regarded as being composed of different kernels  $K_1, K_2, \dots$ , and  $K_n$ .

### 1. The final schedule respects any zero dependence.

Given a dependence  $(a \rightarrow b, \delta, \langle d_1, 0, \dots, 0 \rangle)$ ,  $\mathbf{I} = (i_1, i_2, \dots, i_n)$ , and  $\mathbf{I}' = (i_1 + d_1, i_2, \dots, i_n)$ , we prove that  $f(b, \mathbf{I}') - f(a, \mathbf{I}) \geq \delta$ .

For convenience, denote the sub-vector  $(i_2, \dots, i_n)$  as  $\mathbf{I}_{2n}$ . Then we can represent  $\mathbf{I}$  as  $(i_1, \mathbf{I}_{2n})$ , and  $\mathbf{I}'$  as  $(i_1 + d_1, \mathbf{I}_{2n})$ .

$$\begin{aligned} f(b, (i_1 + d_1, \mathbf{I}_{2n})) - f(a, (i_1, \mathbf{I}_{2n})) &= \underbrace{f(b, (i_1 + d_1, \mathbf{I}_{2n})) - f(b, (i_1 + d_1 - 1, \mathbf{I}_{2n}))}_{(A)} + \\ &\underbrace{f(b, (i_1 + d_1 - 1, \mathbf{I}_{2n})) - f(b, (i_1 + d_1 - 2, \mathbf{I}_{2n}))}_{(B)} + \\ &\dots \\ &\underbrace{f(b, (i_1 + 1, \mathbf{I}_{2n})) - f(b, (i_1, \mathbf{I}_{2n}))}_{(C)} + \\ &\underbrace{f(b, (i_1, \mathbf{I}_{2n})) - f(a, (i_1, \mathbf{I}_{2n}))}_{(D)}. \end{aligned}$$

Let  $x = \min(\text{level}(a), \text{level}(b))$ , and  $y = \max(\text{level}(a), \text{level}(b))$ . Since  $\text{level}(a) \leq y$ , kernel  $K_{\text{level}(a)}$  must enclose  $K_y$  or is  $K_y$  itself. Similarly,  $K_{\text{level}(b)}$  must enclose  $K_y$  or is  $K_y$  itself as



well. Therefore, the initiation intervals  $T_{level(a)}$  and  $T_{level(b)}$  must be greater than or equal to  $T_y$ . Therefore, in the above formula,

$$\begin{aligned} (A) &\geq T_{level(b)} \geq T_y, \\ (B) &\geq T_{level(b)} \geq T_y, \text{ and} \\ &\dots \\ (C) &\geq T_{level(b)} \geq T_y. \end{aligned}$$

So

$$(A) + (B) + \dots + (C) \geq d_1 * T_y$$

Suppose  $\sigma(a, 0) = p_1 * T_1 + q_1$ , and  $\sigma(b, 0) = p_2 * T_1 + q_2$ , where  $0 \leq q_1, q_2 < T_1$ . That is, operation  $a$  is at the modulo cycle  $q_1$  in stage  $p_1$ , and operation  $b$  is at modulo cycle  $q_2$  in stage  $p_2$ . From operation  $a$  to operation  $b$  in the same iteration point, among all the kernels in this course, the smallest kernel must be  $K_n$ , and the biggest must be  $K_x$ . We can use their initiation intervals  $T_n$  and  $T_x$  to estimate the minimum value of  $(D)$ .

If  $\sigma(b, 0) \geq \sigma(a, 0)$ , then  $(D)$  is nonnegative. The value of  $(D)$  is minimum when all the kernels from  $a$  to  $b$  in the same iteration point have the smallest possible value, which is  $T_n$ . Therefore

$$(D) \geq (p_2 - p_1) * T_n + q_2 - q_1.$$

Thus

$$\begin{aligned} (A) + (B) + \dots + (C) + (D) &\geq d_1 * T_y + (p_2 - p_1) * T_n + q_2 - q_1 \\ &\geq \delta \text{ (by Inequality(23)).} \end{aligned}$$

If  $\sigma(b, 0) < \sigma(a, 0)$ , then  $(D)$  is negative. It is minimum when all the kernels from  $a$  to  $b$  in the same iteration point have the biggest possible value, which is  $T_x$ . Thus

$$(D) \geq (p_2 - p_1) * T_x + q_2 - q_1.$$

Consequently,

$$\begin{aligned} (A) + (B) + \dots + (C) + (D) &\geq d_1 * T_y + (p_2 - p_1) * T_x + q_2 - q_1 \\ &\geq \delta \text{ (by Inequality (23)).} \end{aligned}$$

## 2. The final schedule respects any positive dependence.

As we see from the single II case, respecting the positive dependence requires the satisfaction of the two inequalities (32) and (33) regarding the dependent operations  $a$  and  $b$ , and the imaginary operation  $c$ .

Inequality (32) says that the dependence source  $a(I)$  must have finished after the imaginary  $c(I)$  is issued. After shrinking the kernels in the single-II final schedule, their distance is shortened.

---

**Algorithm 1** A Scheduling Algorithm

---

**Require:** the 1-D simplified DDG for loop  $L_1$ , and the hardware resource description.

**Ensure:**  $\sigma(o, 0)$  for every operation  $o$ .

```
1:  $T_{n+1}=0$ 
2: for  $x = n$  to 1 step -1 do
3:   if  $x < n$  then
4:     Move operations from  $OPSETB_x$  to  $OPSETA_x$ , if dependences allow
5:   end if
6:   for  $T_x = \max(RecMII_x, ResMII_x, T_{x+1})$  to  $thresholdII_x$  step 1 do
7:     ModuloSchedule( $L_x, T_x$ )
8:     if Modulo scheduling succeeded then
9:       if  $x < n$  then
10:        Translate the schedule such that  $K_{x+1}$  starts from modulo cycle 0
11:      end if
12:      Increase  $l_x$  to satisfy the first case of the sequential constraints in Inequality ( 24), if necessary
13:      Goto line 17
14:    end if
15:  end for
16:  Return Failure
17: end for
18: Return Success
```

---

However, note that both are from the same  $OPSETA_x$  ( or  $OPSETB_x$ ), and within the same iteration point. The same is true for all the other operation instances between them. Therefore, the minimum distance between them would be

$$f(c, \mathbf{I}) - f(a, \mathbf{I}) = (l_x + 1) * T_x - 1 - p * T_x - q,$$

which is no less than  $\delta - 1$  according to the sequential constraints. Inequality ( 32) still holds, although the minimum distance changes.

Inequality ( 33) says that the dependence sink appears after the imaginary  $c(\mathbf{I})$ . Shrinking would not change such order, as shrinking removes only null cycles, but changes no relative positions between any operation instances. One may refer to the example in Fig. 10 for intuition.

3. Resources constraints are naturally respected, since removing null cycles does not change the relative positions of the operation instances.

### A.3 A Scheduling Algorithm

There are more than one way to construct a 1-D schedule satisfying the scheduling constraints given in Section 6.4. For example, one might formulate the constraints into a linear programming

problem, where all operations from all levels can be scheduled simultaneously. For another example, one might extend the traditional hierarchical reduction to construct the schedule, as we did for our prototype implementation. We describe the algorithm below.

The traditional hierarchical reduction approach [17] is applied to construct the multi-II kernel nest. Assume we choose the outermost loop for scheduling. We schedule the operations starting from the innermost loop to the outermost loop. To each loop, we apply the traditional modulo scheduling.

The algorithm is shown in Algorithm 1, where  $RecMII_x$  and  $ResMII_x$  are the MII determined by the recurrences and resources within the current loop  $L_x$ ,  $thresholdII_x$  is an estimated maximum II for loop  $L_x$ , and  $ModuloSchedule(L_x, T_x)$  calls the traditional modulo scheduling to schedule loop  $L_x$  under the current II,  $T_x$ .

Before scheduling loop  $L_x$ , we do some preprocessing: move every operation in  $OPSETB_x$  to  $OPSETA_x$ , whenever dependences allow. See line 4 in the algorithm. To some extent, this overcomes the limitation of the kernel nesting constraints that requires the operations in  $OPSETB_x$  must be scheduled after the inner loop. This preprocessing can be optional.

In the experiments, the modulo scheduling method we called is slack scheduling [14], which has already been implemented in the ORC compiler. The input to it is always the simplified 1-D DDG for the outermost loop, whatever loop we are currently scheduling. In scheduling the current loop, consider only the 1-D dependences whose source and sink are in the current loop, including its inner loops, and one of them must be at the current level (In computing  $RecMII_x$ , we considered only these dependences as well).

The innermost loop is scheduled in the same way as the original slack scheduling. So let us focus on an outer loop. Let the current loop under scheduling be  $L_x (x < n)$ . Its inner loop  $L_{x+1}$  has been scheduled and the kernel  $K_{x+1}$  has been produced.  $K_{x+1}$  is treated as an atomic operation with the collective resource requirement of all the operations in it, and a latency of  $S_{x+1} * T_x$ , where  $T_x$  is the current II for scheduling the current loop.  $K_{x+1}$  is scheduled with other operations at loop level  $L_x$ . Note that the latency of this atomic operation is different from the other operations: it is dependent on the value of  $T_x$ . Therefore, if the current loop cannot be successfully scheduled under the current value of  $T_x$ , and the next value is tried, the latency of this operation is changed accordingly.

After  $L_x$  is successfully scheduled, we check whether the atomic operation  $K_{x+1}$  is placed at modulo cycle 0. If not, translate the schedule to make it. That is, if the schedule time of the atomic operation is  $p * T_x + q$ , where  $0 \leq q < T_x$ , we add a number  $T_x - q$  to it and the schedule time becomes  $(p+1) * T_x$ , a multiple of  $T_x$ ; the schedule times of all the other operations are added the same number as well. This will finally produced such a kernel nest as shown in Fig. 23. Every kernel is placed from modulo cycle 0. There is no null cycles between the tops of two kernels. Since we assumed the latency of the atomic operation  $K_{x+1}$  is  $S_{x+1} * T_x$  in scheduling loop  $L_x$ , finally, it covers exactly  $S_{x+1}$  number of stages. This arrangement conveniently ensures that the shape of the kernel nest is consistent with that of the generic kernel nest in Fig. 13, without losing any generality.

Below we describe some details in handling the constraints in the modulo scheduling.

1. Modulo property.



For the dependence constraints in Inequality (23), they are conservatively translated into the following:

$$MinDist(o_1, o_2) = \frac{T_x}{T_n} * (\delta - k * T_y) \quad (41)$$

Let us see how the above inequality is derived. First, from the first case of the dependence constraints, we have

$$\underbrace{(p_2 - p_1) * T_x + q_2 - q_1}_{(A)} \geq \underbrace{\delta - k * T_y}_{(B)} + \underbrace{(p_2 - p_1) * (T_x - T_n)}_{(C)} \quad (42)$$

We know  $0 \leq q_1, q_2 < T_x$ , because they are modulo cycles, and because  $K_x$  starts from modulo cycle 0 due to the specific shape of our kernel nest in Fig. 23. So

$$(C) = \left\lfloor \frac{(A)}{T_x} \right\rfloor,$$

which is no bigger than  $\frac{(A)}{T_x}$ .

Therefore, if we conservatively satisfy the following constraint instead:

$$(A) \geq (B) + \frac{(A)}{T_x} * (T_x - T_n),$$

we can definitely satisfy the original constraint in Inequality (42).

The above constraint translates into

$$(A) \geq \frac{T_x}{T_n} * (B).$$

That is,

$$(p_2 - p_1) * T_x + q_2 - q_1 \geq \frac{T_x}{T_n} * (\delta - k * T_y).$$

Therefore, we let

$$MinDist(o_1, o_2) = \frac{T_x}{T_n} * (\delta - k * T_y). \quad (43)$$

For the second case in the dependence constraints in Inequality 23, we have

$$(p_2 - p_1) * T_x + q_2 - q_1 \geq \delta - k * T_y$$

Thus

$$MinDist(o_1, o_2) = \delta - k * T_y \quad (44)$$

Conservatively, we take the maximum of the two values defined in Equations (43) and (44) as the *MinDist*. That explains the definition we show in Equation (41).

3. Resource Constraints.

Modulo scheduling checks the availability of resources when an operation is scheduled to a cycle. The inner loop of the current loop has been modeled as an atomic operation, whose resource usage is the collective resource usage of all the operations in it. Therefore, resource constraints can be respected naturally by modulo scheduling.

4. Sequential Constraints.

For the first case of the sequential constraints in Inequality (24), we do not know the value of  $l_x$ . It is not known until the current loop  $L_x$  is finished scheduling. Therefore, we delay the enforcement of this case of sequential constraints until  $L_x$  is fully scheduled. Then we increase  $l_x$ , if necessary, to assure the constraints are respected. This is always feasible. See line 12 in Algorithm 1.

For the second case of the sequential constraints in Inequality (24), note that  $f_{x+1} * T_x$  is the schedule time of the atomic operation  $K_{x+1}$  at the current kernel  $K_x$ . Thus the constraint can be equivalently translated into the following:

$$MinDist(o, K_{x+1}) = \delta \quad (45)$$

5. Kernel nesting constraints.

As we said before, in scheduling the current loop  $L_x$  under the current interval,  $T_x$ , the atomic operation representing the inner loop is considered to have a latency of  $S_{x+1} * T_x$ ; and after  $L_x$  is scheduled, the schedule is translated such that the atomic operation starts from modulo cycle 0. This ensures the shape of the kernel nesting shown in Fig. 23.

We also needs to make sure that an operation at the current loop level is not scheduled into the inner loop (the stages the atomic operation takes). This translates into the following:

$$MinDist(o, K_{x+1}) = 1 \quad (46)$$

for every operation  $o$  in  $OPSETA_x(x < n)$ .

And

$$MinDist(K_{x+1}, o) = S_{x+1} * T_x \quad (47)$$

for every operation  $o$  in  $OPSETB_x(x < n)$ .

6. A final note for *MinDist*: this is a basic concept used in traditional modulo scheduling (in line 7 of the algorithm). A more complete picture is as follows: at the beginning of modulo scheduling, for every pair of operations  $o_1$  and  $o_2$ , calculate  $MinDist(o_1, o_2)$ . To do so, the initial value of it is by default  $-\infty$ . Then update the value according to the dependence constraints, sequential constraints, and kernel nesting constraints, i.e., according to the Inequalities (41), (45), (46) and (47). If more than one constraint can be applied, take the maximum value among them. We did not show that above for simplicity. During the scheduling process, when an operation is scheduled or ejected, the *MinDist* array is updated in the same principle.

In our experiments, the modulo scheduling method is essentially the same as the original slack scheduling, except for the calculation of *MinDist*.

## A.4 Loop Rewriting

Loop rewriting (or code generation), corresponding to the last step of SSP, final schedule computation. Instead of specifying a function, it directly generates equivalent code in the form of a new  $n$ -dimensional loop nest. The skeleton of this new loop nest is shown in pseudo-code in Fig. 24. Each  $L'_i$  corresponds to the  $L_i$  loop in the original loop nest. Each *for* loop structure is to be replaced by its equivalent in the target assembly language. Loop  $L'_2$  is what we called before an ILES. It is an  $n - 1$ -dimensional loop. The loop body of the innermost loop  $L'_n$  is an *Innermost Loop Pattern*. After this loop is finished in execution, the next innermost loop  $L'_{n-1}$  is tested. If it has other iterations to run, i.e., the index  $i_{n-1} < N_{n-1} - 1$ , then its pipelines is filled, and it continues to run into the innermost loop of this new iteration of it. This code skeleton is essentially an equivalent, simplified, version of that in the literature [26], with the addition of a finalization phase to be complete <sup>7</sup>.

Code generation has several major aspects to take care of: (1) Generate the operations in the right place. This is a mechanical process based on the code skeleton. (2) Before the new loop nest runs, transfers live-in values of the loop variables to rotating registers (from static registers, for example); and after that, transfers live-out values from rotating registers (to static registers, for example). The initialization and finalization phases in the code skeleton are responsible for them, respectively. (3) Based on the register allocation results, assign correct register names to the register operands of the operations. This is a register renaming issue. (4) Set up the control registers such that the new loop nest runs correctly. That is, ineffective operations are automatically disabled in execution; if the current group of iterations is the last group, the  $L'_1$  loop should ends and the control transfers to the epilog. Incorrect setting might lead to wrong results or dead loop. The initialization phase is responsible for that. (5) Handle the special constraints enforced by the target architecture. For example, the IA-64 architecture does not readily support multi-dimensional loop pipelining: it does not have special branch operations for such pipelining as it does for single loop pipelining; the control registers are changed with the concept of single loop pipelining in mind. The first aspect is essentially machine-independent, while the others are. For simplicity, we do not discuss the details of these concerns. Readers are referred to the literature [26].

The code skeleton is illustrated by an example in Fig. 25. In this specific illustration, there are several general points to be noted: first, as we see from Fig. 25(c), the prolog and the epilog are composed of partial copies of  $K_1$  (More specifically, they are composed of  $f_n - 1$  and  $S_1 - f_n$  number of copies, respectively). Second, the OLP is composed of  $S_n$  number of full copies of  $K_1$ . Third, in the entire  $L'_2$ , there are totally  $S_n$  number of outermost loop iterations running as a group. Fourth, each loop  $L'_x$ , except its inner loops, is composed of stages only from  $K_x$ . For example, in Fig. 25, the stages in the innermost loop  $L'_4$  are from  $K_4$  only; the stages in  $L'_3$  but out of the innermost loop, i.e., the code for filling the  $L_3$  pipelines, are from  $K_3$  only. That is why we can shrink a single-II final schedule to be one with multiple IIs. Finally, each  $L_1$  iteration runs sequentially.

---

<sup>7</sup>Compared with that in the literature, this skeleton removes a control variable *initial<sub>i<sub>n</sub></sub>*, which makes it easier to understand, but may lead to larger code size in certain cases.

```

[Initialization]
[Prolog]
L1': for(i1 = 0; i1 < N1; i1++ = Sn) {
    [Outermost loop pattern]
L2':     for(i2 = 0; i2 < N2; i2++) {
L3':         for(i3 = 0; i3 < N3; i3++) {
    ...
Ln':         for(in = 1; in < Nn; in++)
            [Innermost loop pattern]
            } //end Ln'
            if (in-1 < Nn-1 - 1) {
                [Fill Ln-1 Pipelines]
            }
            ...
        } //end L3'
        if (i2 < N2 - 1) {
            [Fill L2 Pipelines]
        }
    } //end L2'
} //end L1'
[Epilog]
[Finalization]

```

Figure 24: Code Generation for An  $n$ -dimensional Loop Nest

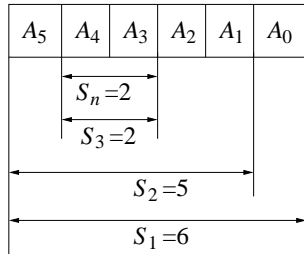


```

L1: for (i1=0; i1 < N1; i1++) {
    OPSETA1
L2:   for (i2=0; i2 < N2; i2++) {
        OPSETA2
L3:     for (i3=0; i3 < N3; i3++) {
L4:       for (i4=0; i4 < N4; i4++) {
            OPSETA4
        } //end L4
    } //end L3
    OPSETB2
  } //end L2
} //end L1

```

(a) A Loop Nest



(b) A 1-D Schedule for Illustration Purpose. Here  $A_x$  represents a stage.

```

[Initialization]
/*Prolog*/
A0
A1 A0
L1': for(i1 = 0; i1 < N1; i1 += 2) {
    /*Outermost loop pattern*/
A5 A4 A3 A2 A1 A0
A5 A4 A3 A2 A1 A0
L2': for(i2 = 0; i2 < N2; i2++) {
L3': for(i3 = 0; i3 < N3; i3++) {
L4': for(i4=1; i4 < N4; i4++) {
    /*Innermost loop pattern*/
A4 A3
A3 A4
  } //end L4'
  if (i3 < N3 - 1) {
    /* Fill L3 Pipelines */
A4 A3
A3 A4
  }
} //end L3'
if (i2 < N2 - 1) {
  /* Fill L2 Pipelines */
A4 A3
A5 A4
A1 A5
A2 A1
A3 A2
}
} //end L2'
} //end L1'
/* Epilog */
A5 A4 A3
A5 A4
A5
[Finalization]

```

(c) The Generated Code

Figure 25: Illustrating the Code Skeleton. The initialization and finalization phases are not shown, as they are machine-dependent.