

# Experiments with auto-parallelizing SPEC2000FP benchmarks

Guansong Zhang, Priya Unnikrishnan, James Ren

IBM Toronto Lab  
Toronto  
ON, L6G 1C7, Canada

**Abstract.** In this paper, we document the experimental work in our attempts to automatically parallelize SPEC2000FP benchmarks for SMP machines. This is not purely a research project. It was implemented within IBM's software laboratory in a commercial compiler infrastructure that implements OpenMP 2.0 specifications in both Fortran and C/C++. From the beginning, our emphasis is on using simple parallelization techniques. We aim to maintain a good trade-off between performance, especially scalability of an application program and its compilation time. Although the parallelization results show relatively low speed up, it is still promising considering the problems associated with explicit parallel programming and the fact that more and more multi-thread and multi-core chips will soon be available even for home computing.

**Keywords:** automatic parallelization, parallelizing compiler, SMT machine, OpenMP, parallel do

## 1 Introduction

Automatic parallelization or auto-parallelization involves automatically identifying and translating serial program code into equivalent parallel code to be executed by multiple threads. Both OpenMP parallelization and auto-parallelization try to exploit the benefits of shared memory, multiprocessor systems. But the main difference between OpenMP and auto-parallelization is that, for OpenMP, the user has complete knowledge of the program behavior and is aware of all the code segments that can benefit from parallelization. For auto-parallelization, however, the challenge is to pick the right parallelizable code from the limited information available within the compiler without any modification to the original source program.

### 1.1 Motivation

It is widely accepted that automatic parallelization is difficult and less efficient than explicit parallel programming. For years, people have been seeking better parallel programming models that can expressively describe the parallelism in an application. The existing models include, just to name a few, HPF[1], MPI[2], OpenMP[3], UPC[4], etc. It is also a known fact that parallel programming is difficult. No matter what kind of parallel programming model is used, creating efficient, scalable parallel code still requires

1 a significant degree of expertise. Parallelization tools can be used as an expert system  
2 to help users to find potential parallelizable code. However, to use the tools effectively,  
3 an understanding of dependence analysis and the application code is still needed.

4 On the other hand, more and more multi-thread and multi-core chips are becoming  
5 available in the market and parallel programming is no longer a privilege available only  
6 for high-performance computing. If not now, in the near future, parallel architecture  
7 will be a common feature in *desktop* and even *laptop computing*. Keeping this growing  
8 trend in mind, auto-parallelization looks like an irresistible option. Auto-parallelization  
9 relieves the user from having to analyze and modify the source program with explicit  
10 compiler directives. The user is shielded from low-level details of iteration space parti-  
11 tioning, data sharing, thread scheduling and synchronization. Auto-parallelization also  
12 saves the user the burden of ensuring correct parallel execution.

13 However, providing compiler support for automatic parallelization is not easy or  
14 straightforward. The biggest critique is that the resultant performance gain is typically  
15 limited, particularly in terms of scalability. In addition, accurate analysis can be re-  
16 stricted by the time constraints a compiler must meet for commercial acceptability,  
17 where in-depth analysis may be sacrificed to allow rapid operation. Taking these factors  
18 into account, we believe that using auto-parallelization is still justified by the following  
19 considerations:

- 20 – For parallel machines with a small number of node processors, scalability is not  
21 a big concern. While auto-parallelization may not give us the speedup that a mas-  
22 sively parallel processing (MPP) application can achieve, it can still be used to take  
23 advantage of the extra silicon available.
- 24 – Even an MPP machine is most likely multi-layered where each node is a tightly  
25 coupled SMP machine, as in a cluster. Explicit message passing parallel programs,  
26 such as MPI, will still benefit from a nested automatic parallel execution.
- 27 – Most of the time-consuming applications are data parallel applications. With the  
28 growing problem size outpacing the improvements in hardware, data parallel appli-  
29 cations can benefit from simple parallelization techniques.
- 30 – As more and more desktop parallel machines are available, users will be willing to  
31 try a parallel approach to solve their computational problems, if only for a small  
32 speedup as long as the effort to achieve that is reasonable.
- 33 – As hardware performance keeps improving, more complicated analysis will be tol-  
34 erable in the compilation process.

## 35 1.2 Organization of the paper

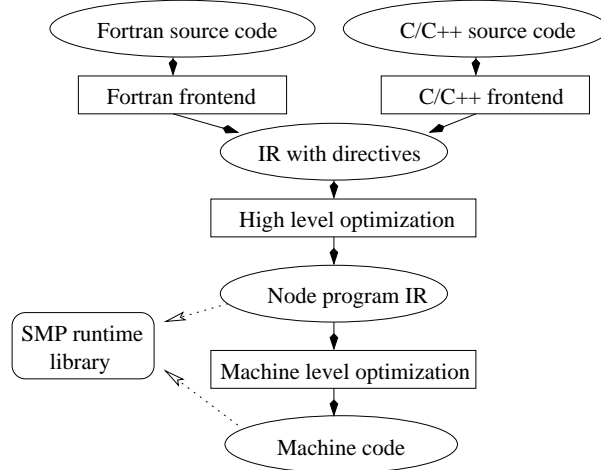
36 The rest of the paper is organized as follows. In Section 2 we introduce our existing  
37 compilation and runtime environment and also briefly describe the structure of the auto-  
38 parallelizer. Section 3 and 4 describes in detail some optimization techniques used by  
39 the parallelizer to enhance the parallel performance of the application code. Section 5  
40 looks at a few challenging parallelization cases that we encountered during the course  
41 of our work. Finally, in Section 6 we summarize our results and future work.

## 2 Parallelizer Structure and Position

The auto-parallelization work done here is based on the available OpenMP compilers, i.e. IBM<sup>®</sup> XL Fortran and VisualAge<sup>®</sup> C/C++. Our auto-parallelizer is essentially a loop parallelizer. For simplicity, the auto-parallelizer avoids complicated parallelization constructs and focuses on identifying and parallelizing expensive loops in the program. The parallelizer simply marks parallelizable loops as OpenMP `parallel do` constructs. The compiler then generates code similar to parallel loops marked by the user. Our auto-parallelizer does not support nested parallelization of loops unlike OpenMP. Auto-parallelization can also be done on an OpenMP program. In this case the auto-parallelizer skips the loop nests with OpenMP constructs and scans for other loops that can potentially benefit from parallelization.

### 2.1 Compilation and Runtime Infrastructure

A typical compiler and runtime structure to support OpenMP is shown in Figure 1 [5]. Figure 2 shows the Compile and Link phases of the compiler. *IR* is the intermediate representation. The parallelizer is delayed until the link phase of the compiler to make maximum utilization of the inter-procedural analysis available during linking. This is shown by the dashed lines in Figure 2



**Fig. 1.** Compile and runtime environment

SPEC2000 CPU benchmark suite was used to evaluate our techniques, as it is one of the best indicators of the performance of the processor, memory and compiler. Our tests focus on Spec2000FP, considering that it is where most of the data parallel processing exists. The hardware system used for testing in this paper is 1.1GHz POWER4<sup>™</sup> system, with 1 to 8 nodes available.

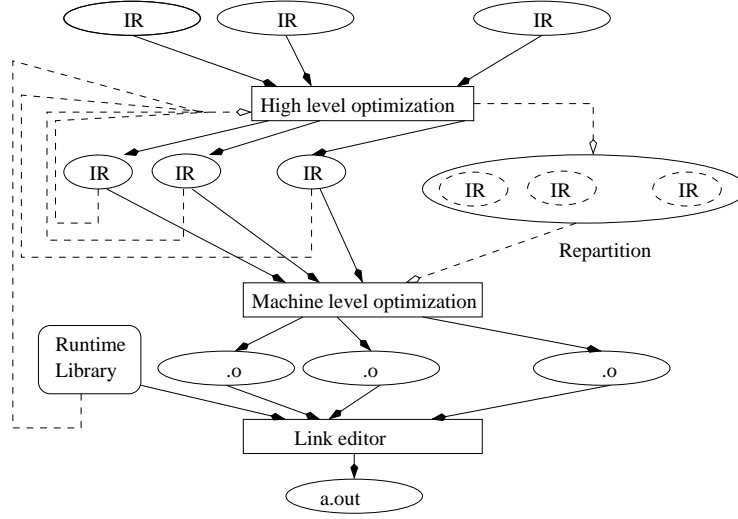


Fig. 2. Compile and link phase optimization

## 2.2 Basic Parallelizer Structure

Dependence analysis is a core part of the parallelizer. Data *dependence vectors*[6] are still our principal tools to analyze and verify the transformation applied on a loop nest. Suppose  $\delta = \{\delta_1 \dots \delta_n\}$  is a hybrid distance/direction vector with the most precise information derivable.

If we have

```

7  L1    DO i1 = 1, U1
8  L2    DO i2 = 1, U2
9          ...
10 Ln    DO in = 1, Un
11 S1    A (f1 (i1, ..., in), ..., fm (i1, ..., in)) = ...
12 S2    ... = A (g1 (i1, ..., in), ..., gm (i1, ..., in))
13 T1    B (u1 (i1, ..., in), ..., um (i1, ..., in)) = ...
14 T2    ... = B (v1 (i1, ..., in), ..., vm (i1, ..., in))
15      END DO
16      ...
17    END DO
18  END DO

```

as a loop nest from  $L_1$  to  $L_n$ , we may have two dependences  $\delta_A$  and  $\delta_B$  characterizing the nested loops, which were caused by  $S_1/S_2$  and  $T_1/T_2$ , where  $f, g, u, v$  are linear functions of the loop induction variables. *Dependence matrix*, denoted as  $\mathcal{D}$ , includes  $\delta_A$  and  $\delta_B$  as two rows.

The Algorithm 1 shows the basic structure of the parallelizer. The parallelizer scans through all the loop nests in a procedure looking for parallelizable loops. Nested par-

---

**Algorithm 1:** Basic loop parallelizer

---

```
begin
  for each loop nest in a procedure do
    for each loop in the nest in the reverse depthfirst order (outer first) do
      if the loop is user parallel then
        ⊥ break
      if the loop is marked sequential, has side-effects etc then
        ⊥ continue
      if the loop has loop carried dependence then
        try splitting the loop to eliminate dependence
        if dependence not eliminated then
          ⊥ continue
      if loop cost is known at compile time then
        if the loop has not enough cost then
          ⊥ break
      else
        ⊥ Insert code for run-time cost estimate
      Mark this loop auto parallel
      break
    end
  end
```

---

1   allelism is avoided by parallelizing only one loop in every nest. The loops in a nest  
2   are scanned from outer to inner as outer loops have a greater benefit from paralleliza-  
3   tion than inner loops. As a first step, loops that are not normalized and loops that are  
4   explicitly marked as sequential or have residuals, side-affects or with loop carried de-  
5   pendences are discarded by the parallelizer at compile time. For loops with loop carried  
6   dependences, the parallelizer tries to split the loop to eliminate the dependences.  
7   The loops are then parallelized independently. This preliminary step eliminates all non-  
8   parallelizable loops. However, naively executing all the parallelizable loops in an ap-  
9   plication may not be a good idea as we will soon discover. Parallelization is a very  
10   expensive operation with high overhead and the compiler has to be very judicious in  
11   identifying loops that can benefit from parallelization.

12   This paper omits discussions about basic techniques of computing dependence ma-  
13   trix, privatizing scalar variables in a loop, finding reduction variables, peeling or split-  
14   ting the loop to eliminate the loop carried dependence, etc. which are all parts of the  
15   preparation phase for the parallelizer. We focus instead on the advanced techniques used  
16   by the parallelizer to intelligently select loops for parallelization and the optimizations  
17   that can be done after identifying a parallelizable loop. Section 3 deals with loop trans-  
18   formation techniques that can enhance the parallel performance of loops selected for  
19   parallelization. Section 4 explores the use of loop cost as an advanced loop selection  
20   technique for parallelization.

### 3 Balancing coarse-grained parallelism and locality

There are many loop transformation techniques that exist today [7]. And almost all of them can be applied to more than one kind of code optimizations. For example, loop interchange is widely used to improve data locality in cache management. Meanwhile, it also can be used to exploit parallelism, such as vectorization, and even coarse-grained parallelism. Unfortunately, the transformations for different optimization purposes do not always work in harmony.

For instance, in a loop nest, maximum spatial cache reuse can be achieved by moving the loop with stride-one access of the references to the innermost position. So a loop permutation algorithm will try to calculate such a loop order while keeping the dependence constraints of the original loop nest. On the other hand, in an automatic parallelizing compiler, one would like to parallelize the outermost loop to reduce the parallel region setup overhead, and leave more code in the inner loop for other optimization opportunities.

There has been significant research to study different approaches to solve this problem. A well known technique is *Unimodular transformation* [8–10], a compound loop transformation algorithm aiming at integrating different loop transformations for a specific goal and target machine, thereby reducing the problem of finding the best transformation combination to finding the best unimodular matrix. Another technique is *Loop selection* [7], which parallelizes all the parallelizable loops, then uses heuristics to select the next one as sequential, expecting to find more parallelizable loops after that. This is easier to implement when compared to unimodular transformation.

Our work is based on the idea of loop selection, but is different from the original one in two aspects: a) given our target is an OpenMP node program, we only need to find one parallelizable loop per nest, to avoid generating code with nested parallelism; b) our heuristics are directly linked to data locality, unlike the original loop selection technique which picks a loop with most “<” directions to parallelize.

#### 3.1 Analyzing model

Using the notation in Section 2, we introduce two theorems.

**Theorem 1 (C-level interchangeable).** *Loop  $L_c$  and  $L_{c-1}$  is interchangeable if and only if  $\forall \delta \in \mathcal{D}$ ,  $\delta \neq (=^{(c-2)} < \dots)$ , where “ $=^{(c-2)}$ ” indicates that there are  $(c-2)$  directions as “=” before the first “<”.*

Theorem 1 and its variant forms can be found in [11] and other literatures. We will not prove it here, but use it directly for our theorem later.

**Theorem 2 (P-level parallelizable).** *Loop  $L_p$  can be parallelized as a parallel do if and only if  $\forall \delta \in \mathcal{D}$ ,  $\delta \neq (=^{(p-1)} < \dots)$ .*

The proof for Theorem 2 is trivial. By the definition of parallel do, it cannot carry any dependence.

Next, we define a *multiply* operation on a dependence vector. Let  $\sigma = \sigma_a \times \sigma_b$ , where  $\forall \sigma_j \in \sigma, \sigma_j = \sigma_{a_j} \times \sigma_{b_j}$ . The multiplication of the two dependence distances is

1 defined as following. It is unknown if one of the distances is unknown. Otherwise, it will  
 2 be the regular multiplication on two integers. If one of the distances is only a direction,  
 3 it is treated as 1 or -1, and after the multiplication, converted back to a direction .

4 Let  $\sigma_p$  and  $\sigma_{p-1}$  be the  $p^{th}$ ,  $(p-1)^{th}$  column vector of dependence matrix  $\mathcal{D}$ , then  
 5 we have

6 **Theorem 3 (Keeping it parallelizable).** *A  $p$ -level parallelizable loop  $L_p$  can be inter-*  
 7 *changed to  $L_{p-1}$  and becomes  $(p-1)$ -level parallelizable if and only if  $\forall \sigma_j \in \sigma$ , where*  
 8  *$\sigma = \sigma_{p-1} \times \sigma_p$ , either*

- 9 -  $\sigma_j = 0$  or
- 10 - the  $j^{th}$  dependence in  $\mathcal{D}$  is not carried by  $L_{p-1}$

### 11 Proof

12 We start by proving the “if” part of the theorem.

13 Suppose  $\forall \sigma_j = 0$ .

14 First, we can assert that  $L_{p-1}$  and  $L_p$  is interchangeable. Otherwise, according to  
 15 Theorem 1 there is a dependence vector  $\delta$  that has the form of  $(=^{(p-2)} < > \dots)$ . This  
 16 will conflict with the fact that all  $\sigma_j$  is zero. Secondly, we prove that the  $L_p$  will become  
 17  $(p-1)$ -level parallelizable after interchanging with  $L_{p-1}$ . Otherwise, from Theorem 2,  
 18 there will a dependence vector  $\delta$ , of the format  $(=^{(p-2)} < \dots)$ , in the dependence matrix  
 19 preventing it from being parallelized after the interchange. Considering the format of the  
 20  $\delta$  before the interchange, given that all  $\sigma_j$  is zero, it should be in the form of  $(=^{(p-1)} < \dots)$ .  
 21 This conflicts with the fact that  $L_p$  is parallelizable according to Theorem 2.

22 If  $\exists \sigma_j \neq 0$ , we let  $\delta$  be the dependence vector on the  $j^{th}$  row of the matrix. It should  
 23 have one of the following formats  $(\dots^{(p-2)} < < \dots)$ ,  $(\dots^{(p-2)} < > \dots)$ ,  $(\dots^{(p-2)} > < \dots)$ ,  
 24  $(\dots^{(p-2)} > > \dots)$ , where  $\dots^{(p-2)}$  is the leading  $(p-2)$  positions. Since  $L_{p-1}$  does  
 25 not carry any dependences as in the given condition,  $\dots^{(p-2)}$  has to be in the format of  
 26  $(= \dots = < \dots)$ , in order for the  $\delta$  to be valid. Therefore,  $L_p$  can be interchanged with  
 27  $L_{p-1}$  and still be kept parallelizable.

28 Similarly, the “only if” part of the theorem can be proved. This is not important in  
 29 our algorithm, so we will omit it here.

30 **End of Proof**

31 **Theorem 4 (Outermost parallelizable).** *Loop  $L_o$  can be parallelized at the outermost*  
 32 *level if and only if  $\sigma = 0$ , where  $\sigma$  is the  $o^{th}$  column vector of dependence matrix  $\mathcal{D}$ .*

33 This is a direct conclusion from Theorem 3. It is actually used in [7] for loop selec-  
 34 tion.

### 35 3.2 Loop permutation with examples

36 Based on our discussion in the previous section, we try to find a permutation favoring  
 37 both data locality and parallelism. To do this, we first assign each loop induction vari-  
 38 able with a weight in terms of memory access. We will favor those having maximum

1 memory spatial reuse. A more comprehensive way of evaluating the weights is called  
2 profitability-based methods, introduced in [7], which has a cache model to estimate  
3 cache misses. Either way, a memory ordering  $O = (L_1, L_2, \dots, L_n)$  of the nested loops  
4 is calculated, where  $L_1$  has the least reuse and  $L_n$  the most.

5 Next we build up a nearby legal permutation in  $P$  by first testing to see if the loop  
6  $L_1$  is legal in the outermost position. If it is legal, it is added to  $P$  and removed from  $O$ .  
7 If it is not legal, the next loop in  $O$  is tested. Once a loop  $L_i$  is positioned, the process  
8 is repeated starting from the beginning of  $O - \{L_i\}$  until  $O$  is empty.  $P$  is considered  
9 as having the best data locality for the loop nest. The algorithm is summarized in [12],

10 Finally, we can adjust the other loops further, based on the following conditions to  
11 decide if  $L_i$  can be moved before  $L_j$ ,

- 12 – The interchange should be legal,
- 13 –  $L_i$  is still parallelizable after moving, and
- 14 – The amount of data locality we are willing to sacrifice.

15 The first two questions are answered by Theorem 3. We will use heuristics including  
16 the induction variable weights calculated previously to control how aggressive the par-  
17 allelizer should be. Apart from estimating the cache miss, a good heuristic could also  
18 include estimations of the loop cost and parallel setup overhead. The topics deserve  
19 some discussions on their own, and will not be included here. *Strip-mining* can also be  
20 considered, if no loop can be moved at all.

21 We illustrate the benefits of using our loop permutation algorithm using a simple  
22 example. Loop permutation for data locality will transform the Fortran code shown  
23 below to  $K$  loop as the outermost and  $I$  loop as the innermost in the nest. This will  
24 lead to the middle  $J$  loop getting parallelized. With our loop permutation algorithm we  
25 can have the  $J$  loop moved to the outermost position, thus striking a balance between  
26 locality and the desired coarse-grained parallelism.

```

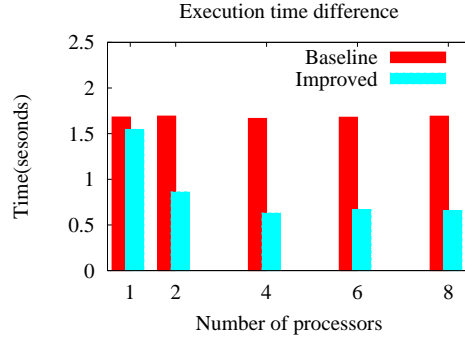
27 DO I = 1, N
28   DO J = 1, N
29     DO K = 1, N
30       A(I, J, K) = A(I, J, K+1) ;
31     END DO
32   END DO
33 END DO

```

34 Figure 3 shows the performance difference of the generated code on the POWER4  
35 system for both the cases discussed.  $N$  is set to 100 and the loop nest is executed 100  
36 times. In the figure, “Baseline” and “Improved” use exactly the same compiler, with  
37 the only difference being that the parallel loop is in the middle for the “Baseline” case  
38 and the outermost for the “Improved” case. We can see that the parallel setup overhead  
39 completely offsets the gain from parallelization in the first case, while the improved  
40 version sees reasonable speedup even for a loop with small computation cost (the loop  
41 body has a single assignment statement).

42 In the SPEC2000FP suite, `mgrid` was negatively affected by this transformation,  
43 further experiments are being carried out to derive better heuristics.





**Fig. 3.** Performance difference

## 4 Restricting parallelization using cost estimation

After a preliminary filtering out of loops unsuitable for parallelization using Algorithm 1, the auto-parallelizer evaluates the cost of a loop to further refine the selection of parallel loops.

### 4.1 Loop Cost Model

The cost of a loop is the approximate execution time of the loop and is given by

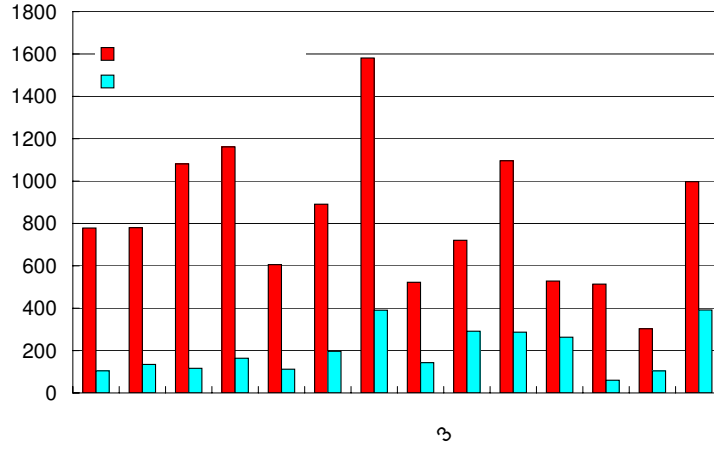
$$\text{LoopCost} = (\text{IterationCount} * \text{ExecutionTimeOfLoopBody})$$

The loop cost algorithm is a recursive algorithm that computes the approximate execution time of the loop including all nested loops inside the current loop. To evaluate loops based on the loop cost, we define a set of *Threshold* values as the basis for comparison. The threshold values are carefully chosen values based on heuristics and take into account the number of processors, overhead arising from the setup required for parallelization, etc.

For loops whose cost can be estimated at compile-time, the parallelizer simply compares the cost against the precomputed threshold value and rejects loops with low costs. However for loops whose cost cannot be determined at compile-time (which is mostly the case), the auto-parallelizer computes a cost expression for the loop in terms of the values known only at runtime. This expression is then evaluated at runtime to determine whether the loop is worth parallelizing. The runtime cost expression is kept as lightweight as possible because of the overhead incurred in computing the expression at runtime. The auto-parallelizer also has a mechanism to limit the number of threads used to parallelize a loop depending on the cost of the loop. Table 1 shows the parallel loops selected at different stages of filtering.

Figure 4 shows the effectiveness of using the Loop Cost Model to restrict parallelization. The results show that careful restriction and appropriate selection of loops

1 for parallelization is extremely crucial for parallel performance. Evaluation of our auto-  
 2 parallelization techniques for Spec2000FP benchmarks indicate that the auto-parallelizer  
 3 is quite precise in selecting high-cost loops for parallelization as shown by Table 2. The  
 4 Spec2000FP benchmarks were analyzed using Profile Directed Feedback techniques  
 5 [13] [14] to identify the high-cost loops in the program. The loops selected by the auto-  
 6 parallelizer are then compared with the list of expensive loops obtained from profile  
 7 directed feedback. From Table 2 we can see that the auto-parallelizer is quite accurate  
 8 within a small error margin. Column 4 of Table 2 indicates the number of loops incor-  
 9 rectly picked by the parallelizer, i.e., loops that have a low cost and hence are not worth  
 10 parallelizing. It is important to keep this number a minimum as selecting the wrong loop  
 11 can adversely impact the parallel execution performance. The zero values in Column 4  
 12 indicate that the parallelizer never picks the wrong loops.



**Fig. 4.** Benchmark performance with and without Loop Cost using two CPU's

## 13 **4.2 Runtime Profiling**

14 In addition to using loop cost for restricting parallel loops, the auto-parallelizer employs  
 15 runtime profiling[15] to further filter out loops at a finer granularity. The runtime pro-  
 16 filing is a more accurate way of measuring the execution time of a chunk of code. The  
 17 execution time of a chunk of a loop executed in parallel (as measured by the runtime  
 18 profiler) is compared with selected *serial* and *parallel* thresholds. If the execution time  
 19 of the chunk is less than the serial threshold, the next chunk of the loop is serialized dis-

Benchmark	#Input Loops	#Loops Processed	#After Stage1	#After Stage2/	#After Stage3
<b>Stage1: Preliminary compile-time filtering</b> <b>Stage2: Compile-time loop cost filtering</b> <b>Stage3: Runtime loop cost filtering</b>					
swim	23	18	9	9	5
wupwise	58	58	6	6	4
mgrid	47	30	11	11	7
applu	180	104	27	27	11
galgel	642	538	280	280	36
lucas	123	123	16	16	4
mesa	777	777	118	118	0
art	83	83	18	18	5
equake	62	64	0	0	0
ammp	283	282	18	18	0
apsi	303	284	92	92	5
sixtrack	1157	1135	347	347	11
fma3d	1648	1635	272	272	33
facerec	202	133	69	69	9

**Table 1.** Stages of Loop Selection for Auto-parallelization of Spec2000 FP Benchmarks

Benchmark	#HighCostLoops from PDF	#Parallelizable HighCostLoops from PDF	#HighCostLoops selected by Parallelizer	#LowCostLoops selected by Parallelizer
swim	8	5	5	0
wupwise	7	4	4	0
mgrid	9	7	7	0
applu	41	11	11	0
galgel	84	49	36	0
lucas	47	4	4	0
mesa	74	3	0	0
art	37	5	5	0
equake	8	0	0	0
ammp	31	0	0	0
apsi	28	11	5	0
sixtrack	44	13	11	0
fma3d	61	33	33	0
facerec	43	25	9	0

**Table 2.** Evaluation of Loop Selection for Auto-parallelization of Spec2000 FP Benchmarks

1 regarding the decision from the runtime loop cost evaluation. The decision to parallelize  
2 or serialize changes dynamically depending on previous executions of the loop.

## 3 **5 Missed parallelization opportunities**

4 As expected, the performance improvements from auto-parallelization are limited (refer  
5 to Section 6 for the actual data). To understand the auto-parallelization results, we com-  
6 pare our results with the SPECOMP benchmarks. For an SMP machine with a small  
7 number of processors, SPECOMP achieves relatively good performance and scalabil-  
8 ity. With this comparison, we hope to understand the reasons for the disparity in the  
9 results between explicit and auto parallelization and also expose opportunities missed  
10 by the auto-parallelizer.

11 Among the 14 SPEC2000FP benchmarks, 10 of them are also included in SPECOMP  
12 test suite. The two versions of the benchmarks were compared on a loop-to-loop basis.  
13 This analysis exposes limitations in the auto-parallelizer and also led to some very inter-  
14 esting observations. We list here some of the cases where the auto-parallelizer failed to  
15 parallelize a loop that was explicitly parallelized in the SPECOMP version. We believe  
16 that improvements to the auto-parallelizer to handle these cases can result in significant  
17 performance improvements for at least some of the benchmarks. To prove this, we try to  
18 manually parallelize loops that were explicitly parallelized in SPECOMP and compare  
19 the performance gain (manual parallelization is not possible in all cases).

### 20 **5.1 Case 1: Loop body contains function calls**

21 The auto-parallelizer fails to parallelize loops that have function calls in the loop body.  
22 Function calls could result in side-effects; the current auto-parallelizer is not capable of  
23 handling such loops. Shown here is a code snip found in *wupwise* where the loop body  
24 has function calls. *wupwise* has four such case, two in *muldeo.f* and two in *muldoe.f*.  
25 Manual parallelization of such loops shows good speedup. To manually parallelize the  
26 loops, array's AUX1 and AUX3 were privatized. The execution time is about 114 sec-  
27 onds with one thread, 61.8 seconds for two threads and 46.5 seconds with four threads,  
28 i.e., a 46% improvement with 2 threads and 59% improvement with 4 threads over using  
29 1 thread. More complicated forms of this case exists in *apsi*, *galgel* and *applu*

```

30      COMPLEX*16      AUX1 (12) , AUX3 (12)
31      . . . .
32
33      DO 100 JKL = 0, N2 * N3 * N4 - 1
34          L = MOD (JKL / (N2 * N3), N4) + 1
35          LP=MOD (L,N4)+1
36          K = MOD (JKL / N2, N3) + 1
37          KP=MOD (K,N3)+1
38          J = MOD (JKL, N2) + 1
39          JP=MOD (J,N2)+1
40          DO 100 I=(MOD (J+K+L,2)+1) , N1, 2
41              IP=MOD (I,N1)+1

```

```

1          CALL GAMMUL(1,0,X(1,(IP+1)/2,J,K,L),AUX1)
2          CALL SU3MUL(U(1,1,1,I,J,K,L),'N',AUX1,AUX3)
3          CALL ZCOPY(12,AUX3,1,RESULT(1,(I+1)/2,J,K,L),1)
4 100      CONTINUE

```

## 5.2 Case 2: Array privatization

Several loops in SPECOMP benchmarks have arrays that are explicitly marked as private, enabling the loops to be parallelized. Array privatization analysis is not yet implemented in our parallelizer preventing it from exploiting these opportunities. The code snip in Case 1 illustrates an example where array privatization is required in order to parallelize the loop. 3 such nested loops can be found in subroutine *rhs* of *applu*. Manual parallelization of such loops resulted in a 12% performance gain for *applu*. Similar situations exist in *galgel* and *apsi*

## 5.3 Case 3: Zero trip loops

Zero trip loops are loops in which the number of iterations calculated from the parameters of the loop is less than 1. There are 8 such loops in *apsi*, one of which is shown below. There exists loop carried dependence on the induction variable L, which the compiler tries to eliminate by rewriting the induction variable in terms of the loop parameters. However, in the case shown here, the parallelizer cannot identify L as an induction variable because of the possibility that the value of NX or NY is zero, thereby preventing the outer loop from getting parallelized. By manually parallelizing the outermost loop of such loops, the performance of *apsi* was improved by 3%.

```

22      DO 30 K=1,NZTOP
23      DO 20 J=1,NY
24      DO 10 I=1,NX
25      L=L+1
26      DCDX(L)=- (UX(L)+UM(K))*DCDX(L)-(VY(L)+VM(K))*DCDY(L)
27      *      +Q(L)
28 10      CONTINUE
29 20      CONTINUE
30 30      CONTINUE

```

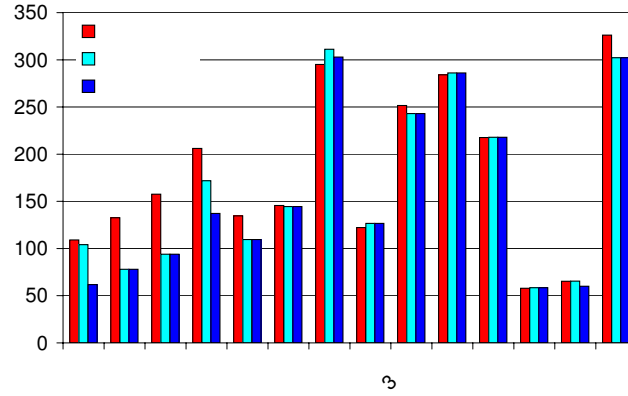
## 6 Summary and future work

We measured the performance of both the sequential and the auto-parallel versions of SPEC2000FP benchmarks. The results are shown in Figure 5<sup>1</sup>. The figure shows three execution times for each benchmark,

- *Sequential* : sequential execution time of the benchmark code compiled by our compiler at the highest optimization level.

<sup>1</sup> The numbers shown here are based on a snap shot of our development compiler.

- 1 – *Parallel* : parallel (2 threads) execution time of the benchmark parallelized by the
- 2 auto-parallelizer for SMP machines.
- 3 – *Improved* : parallel (2 threads) execution time with manual changes to the source
- 4 code as described in Section 5



**Fig. 5.** Performance of auto-parallelization

5 We also emphasize that for simplicity, the modification to individual benchmarks  
6 for the *Improved* results address only one specific problem among the 3 identified in  
7 Section 5. For instance, *apsi* is affected by zero trip loops, function calls in the loop  
8 body as well as array privatization. However, the source modification includes only  
9 source changes for zero trip loops. Among those 10 benchmarks in the SPECOMP test  
10 suite, *swim*, *mgrid*, *applu*, *galgel* and *wupwise* see reasonable speedup, with  
11 improvement up to 40%.

12 The results from auto-parallelization as shown by Figure 5 are not very impres-  
13 sive, especially considering that the hardware resource has actually doubled. However,  
14 we argue that the focus in this paper has been the utilization of simple techniques to  
15 achieve parallelization with minimal impact on the compilation time. This paper has  
16 also shown that much better performance is obtainable using sophisticated techniques  
17 at the cost of increased compilation time. In addition, this paper also presents a simple  
18 loop permutation algorithm to balance data locality and coarse-grained parallelism on  
19 SMP machines.

20 Our auto-parallelizer has several limitations as shown in Section 5. Implementing  
21 array privatization analysis and improving the existing dependence analysis are part of

1 the planned future work. We also plan to fine-tune several heuristics and guidelines  
2 used in the parallelizer. With all these improvements, we hope that the auto-parallelizer  
3 will be able to detect inherent parallelism in the application code on par with explicit  
4 parallelization if not better.

## 5 **7 Trademarks and copyright**

6 IBM, POWER4, and VisualAge are trademarks or registered trademarks of Interna-  
7 tional Business Machines Corporation in the United States, other countries, or both.  
8 Other company, product and service names may be trademarks or service marks of  
9 others.<sup>2</sup>  
10 © Copyright International Business Machines Corporation, 2004. All rights reserved.

## 11 **References**

- 12 1. Charles H. Koebel, David B. Loveman, and Robert S. Schreiber. *The High Performance*  
13 *Fortran Handbook*. MIT Press, 1993.
- 14 2. Peter Pacheco. *Parallel Programming with MPI*. Morgan Kaufmann, 1996.
- 15 3. Rohit Chandra et al. *Parallel programming in OpenMP*. Morgan Kaufmann Publishers,  
16 2001.
- 17 4. Tarek A. El-Ghazawi, William W. Carlson, and Jesse M. Draper. Upc language specification  
18 (v 1.1.1), 2003. <http://upc.gwu.edu>.
- 19 5. Guansong Zhang, Raul Silvera, and Roch Archambault. Structure and algorithm for imple-  
20 menting OpenMP workshare. In *WOMPAT*, Lecture Notes in Computer Science. Springer,  
21 2004.
- 22 6. U. Banerjee. *Dependence Analysis for Supercomputing*. Boston MA: Kluwer, 1988.
- 23 7. Randy Allen and Ken Kennedy. *Optimizing compilers for modern architectures*. Morgan  
24 Kaufmann Publishers, 2002.
- 25 8. Michael E. Wolf and Monica S. Lam. A loop transformation theory and an algorithm to  
26 maximize parallelism. *IEEE Transactions on Parallel and Distributed Systems*, 2(4):452–  
27 471, 1991.
- 28 9. Michael E. Wolf and Monica S. Lam. A data locality optimizing algorithm. *ACM SIG-*  
29 *PLAN'91 Conference on Programming Language Design and Implementation*, 1991.
- 30 10. U. Banerjee. Unimodular transformations of double loops. In *Proc. of the 3rd Workshop on*  
31 *Programming Languages and Compilers for Parallel Computing*, Irvine, CA, August 1990.
- 32 11. Hans Zima and Barbara Chapman. *Supercompilers for Parallel and Vector Computers*.  
33 Addison-wesley, 1990.
- 34 12. K.S.Mckinley, S.Carr, and C.-W Tseng. Improving data locality with loop transformations.  
35 *ACM Trans. on Programming Language and Systems*, 18(4), 1996.
- 36 13. Robert Cohn and P. Geoffrey Lowney. Feedback directed optimization in Compaq's compi-  
37 lation tools for Alpha. 2nd ACM Workshop on Feedback-Directed Optimization, 1999.
- 38 14. W. J. Schmidt et. al. Profile-directed restructuring of operating system code. In *IBM Systems*  
39 *Journal*, 37(2), 1998.
- 40 15. Sagnik Nandy, Xiaofeng Gao, and Jeanne Ferrante. TFP: Time-sensitive, Flow-specific Pro-  
41 filing at Runtime. In *LCPC 2003*, 16th Workshop on Languages and Compilers for Parallel  
42 Computing, 2003.

---

<sup>2</sup> The opinions expressed in this paper are those of the authors and not necessarily of IBM.