

Minimum Register Instruction Sequencing to Reduce Register Spills in Out-of-Order Issue Superscalar Architectures

R. Govindarajan, *Senior Member, IEEE*, Hongbo Yang, *Student Member, IEEE*,
José Nelson Amaral, *Senior Member, IEEE*, Chihong Zhang, *Member, IEEE*, and
Guang R. Gao, *Senior Member, IEEE*

Abstract—In this paper, we address the problem of generating an optimal instruction sequence S for a Directed Acyclic Graph (DAG), where S is optimal in terms of the number of registers used. We call this the Minimum Register Instruction Sequence (MRIS) problem. The motivation for revisiting the MRIS problem stems from several modern architecture innovations/requirements that has put the instruction sequencing problem in a new context. We develop an efficient heuristic solution for the MRIS problem. This solution is based on the notion of *instruction lineage*—a set of instructions that can definitely share a single register. The formation of lineages exploits the structure of the dependence graph to facilitate the sharing of registers not only among instructions within a lineage, but also across lineages. Our efficient heuristics to “fuse” lineages further reduce the register requirement. This reduced register requirement results in generating a code sequence with fewer register spills. We have implemented our solution in the MIPSpro production compiler and measured its performance on the SPEC95 floating point benchmark suite. Our experimental results demonstrate that the proposed instruction sequencing method significantly reduces the number of spill loads and stores inserted in the code, by more than 50 percent in each of the benchmarks. Our approach reduces the average number of dynamic loads and stores executed by 10.4 percent and 3.7 percent, respectively. Further, our approach improves the execution time of the benchmarks on an average by 3.2 percent. In order to evaluate how efficiently our heuristics find a near-optimal solution to the MRIS problem, we develop an elegant integer linear programming formulation for the MRIS problem. Using a commercial integer linear programming solver, we obtain the optimal solution for the MRIS problem. Comparing the optimal solution from the integer linear programming tool with our heuristic solution reveals that, in a very large majority (99.2 percent) of the cases, our heuristic solution is optimal. For this experiment, we used a set of 675 dependence graphs representing basic blocks extracted from scientific benchmark programs.

Index Terms—Compiler optimization, code sequence optimization, register allocation, instruction scheduling, code generation, superscalar architectures, instruction level parallelism.

1 INTRODUCTION

IN this paper, we revisit the *optimal code generation* problem [1], [9], also known as the *evaluation-order determination* problem [42]: the problem of generating an instruction sequence from a data dependence graph (DDG). In particular, we are interested in generating an instruction sequence S that uses the minimum number of registers. We define the *Minimum Register Instruction Sequence* (MRIS) problem as:

Problem Statement 1. *Given a data dependence graph G , derive an instruction sequence S for G that is optimal in the sense that its register requirement is minimum.*

- R. Govindarajan is with the Supercomputer Education & Research Centre and Department of Computer Science & Automation, Indian Institute of Science, Bangalore, India.
- H. Yang and G.R. Gao are with the Electrical and Computer Engineering Department, University of Delaware, Newark, DE 19716.
- J.N. Amaral is with the Department of Computing Science, University of Alberta, Edmonton, AB, T6G-1R4, Canada.
E-mail: amaral@cs.ualberta.ca.
- C. Zhang is with Qualcomm Inc., 5775 Morehouse Dr., San Diego, CA 92121.

Manuscript received 11 June 2001; revised 2 Apr. 2002; accepted 3 Apr. 2002.
For information on obtaining reprints of this article, please send e-mail to: tc@computer.org, and reference IEEECS Log Number 114328.

Our study of the MRIS problem is motivated by the challenges faced by modern processor architectures. For example, a number of modern superscalar processors support out-of-order instruction issue and execution [37]. Out-of-order (o-o-o) instruction issue is facilitated by the runtime scheduling hardware and by the register renaming mechanism in superscalar architectures. An o-o-o processor has more physical (i.e., not logical or architected) registers at its disposal for register renaming at runtime that are not visible to the compiler. Due to these hardware mechanisms, o-o-o issue processors have the capability to uncover the instruction level parallelism obscured by anti and output-dependences (together known as *false dependences* or *name dependences*). Hence, it is important in these processors to reduce the number of register spills, even at the expense of not exposing instruction level parallelism [36], [40]. Reducing register spills reduces the number of loads and stores executed, which in turn is important:

- from a performance viewpoint in architectures that either have a small cache or a large cache miss penalty;
- from a memory bandwidth usage viewpoint;

- from an instruction-level parallelism viewpoint as the elimination of some of the spill instructions frees instruction slots to issue other useful instructions;
- from a power dissipation viewpoint, as load and store instructions contribute to a significant portion of the power consumed.

Another argument as to why the MRIS problem is relevant is the code generation for threads in fine-grain multithreaded architectures. In such architectures, it is often important to minimize the number of registers used in a thread in order to reduce the cost of a thread context switch [14].

The MRIS problem is related to, but different from, the conventional instruction scheduling [1], [16], [17], [28], [41] and the register allocation [1], [7], [8], [10], [11], [15], [28], [30], [32], [39] problems. The main objective of a traditional instruction scheduling method is to minimize the execution time (or length) of the schedule. Hence, an instruction scheduling method must take into account the execution latencies of each instruction in the DDG. In contrast, the latency of the instructions in the DDG and the availability of functional unit resources are not a part of the MRIS problem formulation. To highlight this difference, throughout this paper, we use the term “*instruction schedule*” to refer to the former and “*instruction sequence*” to refer to the latter.

The MRIS problem is closely related to the *optimal code generation (OCG) problem* [1], [9], [34] or the *evaluation-order determination problem* [42]. For the case in which the dependence graph is a tree, an algorithm that produces an optimal sequence (in terms of code length) exists. For a general DAG, the problem is known to be NP-Complete since 1975 [34]. An important difference between these traditional code generation approaches (OCG and evaluation order determination) and our MRIS problem is that the former emphasize reducing the schedule length, while the latter focuses on minimizing the number of registers used.

In this paper, we present a simple and efficient heuristic method to address the MRIS problem. The proposed method is based on the following:

- *Instruction lineage formation*: The concept of an *instruction lineage* evolves from the notion of an *instruction chain* [19] which allows the sharing of a register among instructions along a (flow) dependence chain in a DDG. In other words, a lineage is a set of nodes in the DDG that can definitely reuse the same destination register. Instruction lineages model the DDG register requirement more accurately than instruction chains (see Section 6).
- *Lineage Interference Graph*: The notion of a lineage interference graph captures the *definite overlap* relation between the live ranges of lineages even before the instructions are scheduled and it is used to facilitate sharing of registers across lineages.

We have implemented our heuristic approach in the SGI MIPSpro compiler suite and evaluated its performance on the SPEC95 floating-point benchmarks. For comparison, we also measured the performance of a baseline version of the compiler that executes traditional compiler optimizations, but that does not optimize the instruction sequences with the goal of reducing register pressure. We also compare our

approach with an Optimized version of the MIPSpro compiler that performs integrated instruction scheduling and register allocation. As the emphasis of our work is on sequencing the instructions to reduce the register requirements and the number of spill instructions executed, we measure the number of static as well as dynamic loads and stores in each benchmark under different versions of the compiler. We also compare the execution time of the benchmarks. Our experimental results are summarized as follows:

- When compared to the baseline version of the MIPSpro compiler, our heuristic approach significantly reduces the number of spill instructions inserted in the (static) code, on an average by 63 percent, and by more than 50 percent in each of the benchmarks. The heuristic approach also drastically reduces the number of basic blocks that require spills.
- Our heuristic method reduces the number of dynamic load and store instructions executed by as much as 20.9 percent and 11.2 percent, respectively, and on average by 10.4 percent and 3.7 percent, in comparison to the baseline version of the compiler.
- Our heuristic method also marginally improves the execution time, on average by 3.2 percent. In one of the applications, the performance improvement is as high as 14.2 percent.
- Last, even compared to the optimized version, our heuristic approach performs better in terms of all the above parameters, although the percentage improvement is relatively low.

The baseline version of the compiler performs the same optimizations as in the optimized compiler, including the same instruction scheduling and register allocation algorithms. In the optimized compiler, when the local register allocation requires spilling, the local instruction scheduler is called again, using a more accurate estimate for the register pressure. This second attempt at instruction scheduling optimization, after the first attempt at register allocation, is disabled in the baseline compiler.

In order to assess the optimality and efficiency of the proposed heuristic, we formulate the MRIS problem as an integer linear programming problem and solve it using a commercial integer linear programming solver. This implementation produces the optimal solution for the MRIS problem for DDGs with a small number of nodes. Comparing the optimal solutions with the heuristic ones reveals that our heuristic solution was optimal for a very large majority (99.2 percent) of the cases. For this experiment, we used a set of 675 (small sized) DDGs representing basic blocks in scientific benchmarks programs.

The rest of the paper is organized as follows: In the following section, we motivate the MRIS problem with the help of an example. In Section 3, we present our heuristic solution for lineage formation and a sequencing method to construct a near-optimal minimum register instruction sequence.¹ Section 4 deals with the formulation of the MRIS

1. We use the term *near-optimal* to indicate that in our empirical investigation, our heuristic algorithm compared well with known optimal solutions. We make no claims about near-optimality from an approximation theory standpoint.

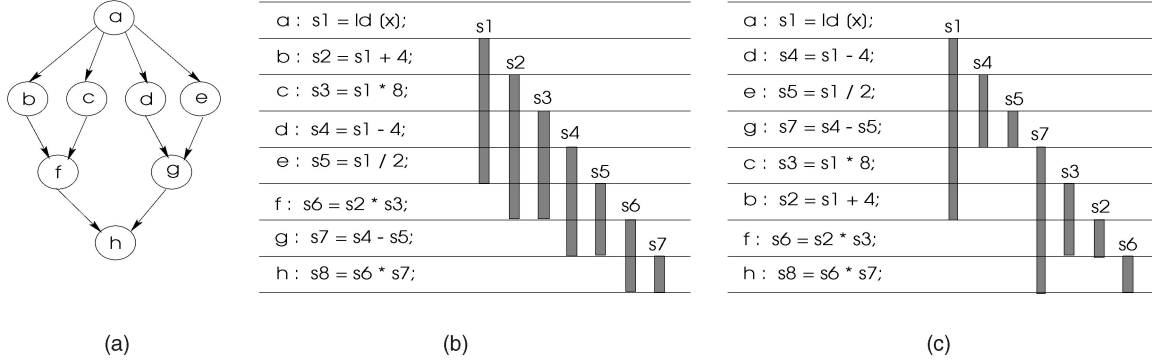


Fig. 1. Motivating example. (a) DDG. (b) Instruction sequence 1. (c) Instruction sequence 2.

problem as an integer linear programming problem. We report static and dynamic performance measures of our approach in Section 5. Related work and conclusions are presented in Section 6 and 7.

2 MOTIVATION AND OVERVIEW

In this section, we first use an example to motivate the MRIS problem. Later, we illustrate our heuristic approach using the same example.

2.1 Motivating Example

Consider the computation represented by the data dependence graph (DDG)² shown in Fig. 1a. Two possible instruction sequences for this DDG are also shown in the figure, along with the live ranges of variables $s1 - s7$ (for simplicity, we assume, in this example, that all the variables are dead at the end of the basic block). For the instruction ordering shown in Fig. 1b, we have four variables simultaneously live in statements e and f , therefore, four registers are required. However, for the sequencing shown in Fig. 1c, only three variables are simultaneously live and, therefore, we may use only three registers. In this particular example, the minimum register requirement is three. Hence, the sequence shown in Fig. 1c is one of the minimum register sequences.

The MRIS problem can be stated as follows: Given a set of instructions and the data dependences among them, construct an instruction sequence that requires the minimum number of registers. The input for the MRIS problem is a Data Dependence Graph (DDG) where the nodes represent instructions and the directed edges, also referred to as flow arcs, represent data dependences.³ The edges of the DDG impose a partial order among the instructions. In this paper, we restrict our attention to acyclic DDGs and, hence, do not consider any loop-carried dependences.

We say that multiple instructions share a single register if they use the same register as the destination for the values

that they produce. We need to identify which nodes in the DDG can share the same register in all legal sequence. Although a complete answer to this question is hard to determine, the data dependences in the DDG provide a partial answer. For instance, in the DDG of Fig. 2a, since there is a data dependence from node b to node f and there are no other nodes that use the value produced by node b , we can definitely say that, in any legal sequence, the register associated with node b can be shared by node f . Similarly, nodes e and g can share the same register. Next, can nodes f and g share the same register? The answer is no because, in any legal sequence, the values produced by these instructions must be live simultaneously so that the computation in node h can take place.

Another interesting question is whether nodes c and d can share the same register. The data dependences in the DDG neither require their live ranges to definitely overlap (as in the case of nodes f and g) nor imply that they definitely will not overlap (as in the case of nodes b and f). Hence, to obtain a minimum register instruction sequence, one must order the nodes in such a way that the live ranges of nodes c and d do not overlap and, hence, they can share the same register. In the following subsection, we outline our approach, which uses efficient heuristics to arrive at a near-optimal solution to the MRIS problem.

2.2 Overview of the Lineage-Based Algorithm

Our solution to the MRIS problem uses the notion of *instruction lineages*. If S_i is a node in a DDG and S_i has one or more descendants, then S_i produces a value and must be assigned a register.⁴ If we have a sequence of instructions in the DDG S_1, S_2, \dots, S_n , where S_2 is the descendant of S_1 , S_3 is the descendant of S_2 , etc., then we can form a *lineage* of these instructions in such a way that all the instructions in the lineage share the same register. That is, the register assigned to S_1 is passed on to S_2 (S_1 's *heir*), which is passed on to S_3 , and so on. Due to the data dependences between pairs of instructions in the lineage, any legal sequence will order the instructions as S_1, S_2, \dots, S_n . Hence, if the live

2. Since our focus in this paper is on generating an instruction sequence that reduces register pressure, we consider only flow dependences in our DDG. Other dependences due to memory (such as store-load dependences), while important from a scheduling viewpoint, do not influence register allocation and hence need not be considered for computing the register requirements.

3. Although we present the MRIS formulation for the DDG of a basic block, our method is also applicable to superblocks [23].

4. A node of the DDG without descendants must either be a store instruction or produce a value that is live-out of the basic block. For simplicity, we do not consider live-out registers in the presentation of our solution. A simple extension of our solution would insert a dummy sink node to capture live-out registers, as in [27].

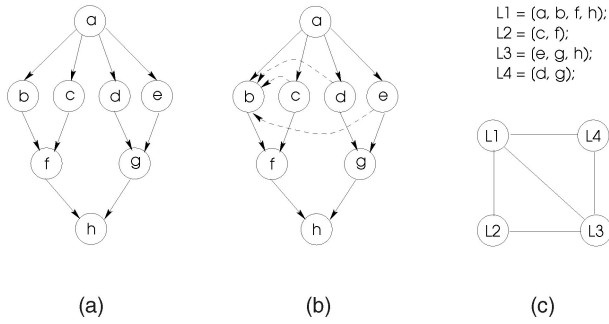


Fig. 2. Data dependence graph for the motivating example. (a) Original DDG. (b) Augmented DDG. (c) Lineage interference graph.

range of the variable defined by S_1 ends at S_2 , then S_2 can definitely reuse the same register allocated to S_1 .

What if S_1 has more than one descendant? In order for S_2 to use the same register that S_1 used, we must ensure that the other descendants of S_1 are sequenced before S_2 . Thus, the selection of one of the descendants of S_1 to be the *heir* of the register creates sequencing constraints among the descendants of S_1 . Such sequencing constraints are explicitly represented in the *augmented DDG* by means of directed *sequencing edges* from each descendant node to the selected heir. For instance, the DDG for the motivating example is shown in Fig. 2a. The formation of the lineage $L1 = \{a, b, f, h\}$ creates scheduling constraints among the descendants of a . These sequencing edges are shown as dotted arcs in Fig. 2b. In Section 3, we introduce a simple but efficient heuristic to select an heir. We will also show that the introduction of sequencing edges does not introduce cycles in the DDG.

It is clear that all the nodes in a lineage share the same register. But, can two lineages share the same register? To determine the interference between two lineages, we have to determine whether the live ranges of the lineages overlap in all legal instruction sequences. The live range of an instruction lineage is the concatenation of the live ranges of all the values defined by the instructions in the lineage. If the live ranges of two lineages overlap in all legal sequences, then the lineages cannot share the same register. However, if they do not overlap in at least one of the legal sequences, then we may be able to sequence the lineages in such a way that they share a register. Based on the overlap relation, we construct a Lineage Interference Graph (LIG). Fig. 2c shows the LIG for our example. This lineage interference graph is colored using traditional graph coloring algorithms to compute the number of registers required for the DDG [10], [11]. We refer to this number as the *Heuristic Register Bound (HRB)*. Once we color the lineage interference graph, we apply a sequencing method that uses this coloring as a guideline to generate an instruction sequence that uses the minimum number of registers. Our heuristic algorithm produces a near-optimal solution for the MRIS problem.

The MRIS formulation optimizes an instruction sequence based only on the number of registers used in the sequence. It does not take into consideration the size of the instruction window in a superscalar processor. Therefore, an optimal solution to the MRIS problem may result in a sequence that

minimizes the register pressure, but leads to a suboptimal schedule at runtime because instructions that could execute in parallel are not in the same instruction window.

3 HEURISTIC APPROACH TO THE MRIS PROBLEM

In this section, we present our heuristic approach to find a good approximation to a minimum register sequence for an acyclic DDG. First, we formally introduce the concept of lineage and describe our lineage formation algorithm. Next, we establish a condition for lineage overlapping and introduce the lineage interference graph (LIG) in Section 3.2. This section also deals with a method to compute the HRB. In Section 3.3, we introduce the concept of lineage fusion to improve register sharing. Finally, we describe the sequencing algorithm for instruction sequence generation based on the coloring of the LIG in Section 3.4.

3.1 Lineage Formation

We use the notion of *instruction lineages* to identify sequences of operations that can share registers.

Definition 3.1. An *instruction lineage* $L_v = \{v_1, \dots, v_n\}$ is a set of nodes such that there exist flow arcs $(v_1, v_2), (v_2, v_3), \dots, (v_{n-1}, v_n)$ in the DDG. We say that v_2 is the *heir* of v_1 in lineage L_v , v_3 is the heir of v_2 in L_v , and so on.

Although a node can have many immediate descendants in the DDG, only one of these descendants can be the node's heir. The heir of a node v_i is a node v_j that inherits the destination register of v_i to store its own result. As a consequence, the formation of a lineage imposes sequencing constraints among the descendants of node v_i : The chosen heir v_j must be the last among all the descendants of v_i to be executed because v_j is the *last use* of the value produced by v_i .

The last node v_n in a lineage is the last one to use the value in the register (defined by v_{n-1}) assigned to that lineage. The last node v_n might belong to another lineage or might be a store instruction that does not need a register. Therefore, the last node does not use the same register as the other nodes in the lineage. We emphasize this property of the last node by denoting a lineage with semi-open intervals as $[v_1, v_2, \dots, v_{n-1}, v_n)$. Thus, in the DDG of Fig. 2, the four lineages that cover all the def-last_use relations are $[a, b, f, h)$, $[c, f)$, $[e, g, h)$, and $[d, g)$.

Our lineage formation algorithm attempts to form as long a lineage as possible, in the sense that, if $\{v_1, v_2, \dots, v_n\}$ is a lineage, then either v_n is a node with no descendants or v_n is already associated with some other lineage. Because the heir is always the last descendant to be executed in an acyclic sequencing, the live range of the nodes in a lineage will definitely not overlap with each other and, hence, all the nodes in the lineage can share the same register. In order to ensure that the heir is the last use of the value produced by its parent node, we introduce *sequencing edges* in the DDG from each descendant of a node to the chosen heir, as shown in Fig. 2b. A sequencing edge from node v_i to v_j imposes the constraint that node v_i must be listed before node v_j is listed. This constraint implies that all nodes that can reach v_i must be listed before any node that can be reached from v_j is listed.

If the introduction of sequencing edges were to make the graph cyclic, then it would be impossible to obtain a sequence for the instructions represented in the DDG. Hence, some care should be taken in the selection of an heir. During the formation of instruction lineages, we use a simple height priority to choose the heir of each node. The height of a node is defined as follows:

$$ht(u) = \begin{cases} 1 & \text{if } u \text{ has no descendants} \\ 1 + \max_{v \in D(u)} (ht(v)) & \text{otherwise,} \end{cases}$$

where $D(u)$ is the set of all immediate descendants of u . In the DDG of Fig. 2a in Section 2, the heights of the nodes are:

$$\begin{aligned} ht(a) &= 4; & ht(b) &= 3; & ht(c) &= 3; & ht(d) &= 3; \\ ht(e) &= 3; & ht(f) &= 2; & ht(g) &= 2; & ht(h) &= 1; \end{aligned}$$

During the lineage formation, if a node v_i has multiple descendants, then we choose a descendant node with the smallest height to be the heir of v_i . If multiple descendants have the same lowest height, then the tie is broken arbitrarily. In order to ensure that cycles are not introduced in the lineage formation process, we recompute the height of each node after introducing sequencing edges between the descendants of a node.

Each flow arc (u, v) in a DDG corresponds to a true flow dependence and, hence, to a definition-use (def-use) relationship between the nodes u and v . Hence, each dependence edge is associated with a register. Later, we will assign registers to lineages of nodes of the DDG. Therefore, it is important that the live range of each node of the DDG (except for sink nodes) be associated with exactly one lineage. The lineages are formed by the arcs between nodes and their chosen heirs, i.e., each arc that is part of a lineage is a *def-last_use* arc in the DDG. By using a greedy algorithm to form lineages, our algorithm is conservative in the number of lineages formed, thus reducing the size of the *lineage interference graph* and the complexity of coloring that graph.

The Lineage Formation algorithm is essentially a greedy depth-first graph traversal algorithm that identifies an heir for each node using the height priority. If a node has multiple descendants, sequencing edges are introduced after the heir is selected and the heights of all nodes in the DDG are recomputed. The detailed algorithm is presented in Fig. 3. The application of the lineage formation algorithm to the DDG of Fig. 2a results in four lineages, as shown in Fig. 2c. We refer to the DDG with the additional sequencing edges as the *augmented DDG* (refer to Fig. 2b). Notice that, in Steps 8 and 10 of the algorithm, we distinguish flow edges from sequencing edges. That is, if there is a sequencing edge from node v_i to node v_j , v_j is not considered a descendant of v_i for the purpose of lineage formation. Only arcs that represent data dependences are considered for lineage formation.

Next, we show that the introduction of sequencing edges does not introduce any cycle. Formally,

Lemma 3.1. *The introduction of sequencing edges during lineage formation does not introduce any cycle in the augmented DDG.*

Proof. Since only the lowest descendant of a node is chosen as the heir, all sequencing edges inserted will be from nodes with higher heights to nodes with lower heights. Also, if the lowest descendant is already in a different

```

LINEAGEFORMATION( $V, E$ )
1.  mark all nodes in the DDG as not in any lineage
2.  compute the height of every node in the DDG
3.  while there is a node not in any lineage do
4.      recompute height  $\leftarrow$  false
5.       $v_i \leftarrow$  highest node not in lineage
6.      start a new lineage containing  $v_i$ 
7.      mark  $v_i$  as in a lineage
8.      while  $v_i$  has a descendant do
9.           $v_j \leftarrow$  lowest descendant of  $v_i$ 
10.         if  $v_i$  has multiple descendants
11.             recompute height  $\leftarrow$  true
12.             for each descendant  $v_k \neq v_j$  of  $v_i$  do
13.                 add sequencing edge from  $v_k$  to  $v_j$ 
14.             endfor
15.         endif
16.         add  $v_j$  to lineage
17.         if  $v_j$  is already marked as in a lineage
18.             // end lineage with  $v_j$  as the last node
19.             break;
20.         endif
21.         mark  $v_j$  as in a lineage
22.          $v_i \leftarrow v_j$ 
23.     end while
24.     if recompute height = true
25.         recompute the height of every node in the DDG
26.     endif
27. end while

```

Fig. 3. Lineage Formation Algorithm.

lineage, then the current lineage ends. Furthermore, the heights of all the nodes in the graph are recomputed (in Steps 24 and 25) at the start of each new lineage.⁵ Thus, in the augmented DDG, there can be no path from a node with a lower height to one with a higher height. Therefore, no cycle is formed. \square

3.2 Lineage Interference Graph and Heuristic Register Bound

In this section, we discuss how to determine whether the live ranges of two lineages always overlap and how to compute the heuristic register bound.

3.2.1 Overlapping of Live Ranges

In order to determine whether two lineages can share a register, we need to verify if the live ranges of these lineages overlap. To define the live range of a lineage, we use the fact that each instruction has a unique position t in the sequence of instructions.

Definition 3.2. *If v_1 , the first instruction of a lineage $L = [v_1, v_2, \dots, v_n]$ is in position t_i and the last instruction v_n of L is in position t_j in an instruction sequence, then the **live range** of the lineage L starts at t_{i+1} and ends at t_j .*

Our goal is to find an instruction sequence that requires the minimum number of registers over all legal instruction sequences. Therefore, we must identify lineages that can

5. Although one may think that the recomputation of the heights of all nodes is needed every time after choosing a heir in a lineage, the heights of descendants of the chosen heir will not be affected by the sequencing edges. Hence, the recomputation of heights is performed, if necessary, only when a new lineage is started.

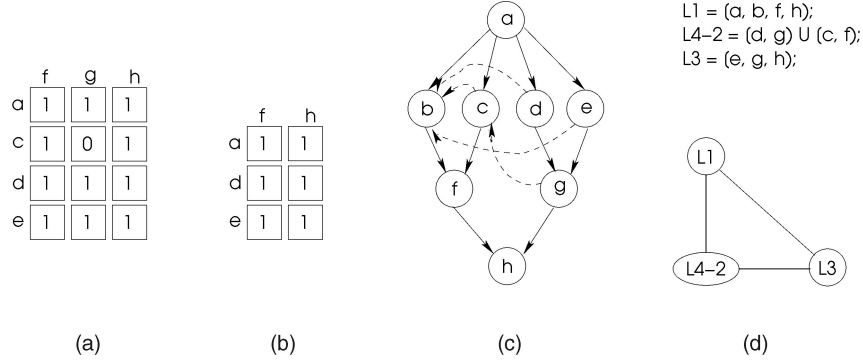


Fig. 4. Reach relation. DDG and LIG after lineage fusion for motivating example. (a) Reachability before fusion. (b) Reachability after fusion. (c) DDG after fusion. (d) LIG after lineage fusion.

share the same register, i.e., lineages with nonoverlapping live ranges. Moreover, we must identify such lineages before a complete instruction sequence is produced. While we are identifying these lineages, the instructions within each lineage are totally ordered, but instructions from multiple lineages will be interleaved in a yet-to-be-determined way to form the final instruction sequence. Thus, we say that the position of instructions in the sequence and, as a consequence, the live ranges of different lineages are *floating*. The live range of a lineage is always contiguous, irrespective of interleavings of instructions from multiple lineages. Thus, once the first instruction in a lineage is listed, its live range is active until the last instruction of the lineage is listed in the sequence.

In order to determine whether the live ranges of two lineages must necessarily overlap, we define a condition based on the existence of paths between the lineages. First, we define two set of nodes: S is the set of nodes that start lineages and E is the set of nodes that end lineages. Next, we define the *reach* relation R .

Definition 3.3. The reach relation $R : S \rightarrow E$ maps S to E . For all $v_a \in S$ and $v_b \in E$, node v_a reaches node v_b , $R(v_a, v_b) = 1$, if there is a path in the augmented DDG from v_a to v_b . Otherwise, $R(v_a, v_b) = 0$.

The reach relation is used to determine whether the live ranges of two lineages must necessarily overlap in *all* legal instruction sequences for the augmented DDG.

Definition 3.4. Let t_{s_i} denote the position of instruction s_i in an instruction sequence. Two lineages, $L_u = [u_1, u_2, \dots, u_m]$ and $L_v = [v_1, v_2, \dots, v_n]$, **overlap** in an instruction sequence if $t_{u_1} < t_{v_1} < t_{u_m}$ or $t_{v_1} < t_{u_1} < t_{v_n}$.

Definition 3.5. Two lineages, $L_u = [u_1, u_2, \dots, u_m]$ and $L_v = [v_1, v_2, \dots, v_n]$, **definitely overlap** if they overlap in all possible instruction sequences.

Theorem 3.1. The live ranges of two lineages, $L_u = [u_1, u_2, \dots, u_m]$ and $L_v = [v_1, v_2, \dots, v_n]$, **definitely overlap** if u_1 reaches v_n and v_1 reaches u_m .

Proof. Since there exists a directed path (consisting of flow and sequencing edges) from u_1 to v_n , u_1 must be listed before v_n in any instruction sequence, i.e., $t_{u_1} < t_{v_n}$. Similarly, since there exists a directed path from v_1 to u_m , $t_{v_1} < t_{u_m}$. Further, since L_u is a lineage, $t_{u_1} < t_{u_m}$ and,

similarly, $t_{v_1} < t_{v_n}$. If u_1 is listed before v_1 , $t_{u_1} < t_{v_1}$, we obtain that $t_{u_1} < t_{v_1} < t_{u_m}$ and the lineages overlap. Likewise, if v_1 is listed before u_1 , $t_{v_1} < t_{u_1}$, we obtain $t_{v_1} < t_{u_1} < t_{v_n}$ and, in this case also, the two lineages overlap. Thus, we have proven that, whether L_u starts before L_v , or L_v starts before L_u , the two lineages overlap. Hence, the live ranges of L_u and L_v definitely overlap. \square

Consider the lineages $L_1 = [a, b, f, h]$ and $L_4 = [d, g]$ in our motivating example in Fig. 2b. Node a can reach node g through the path $a \rightarrow d \rightarrow g$ and node d can reach node h through the path $d \rightarrow g \rightarrow h$. Therefore, the live ranges of these two lineages must overlap. Similarly, the live range of lineage L_2 must overlap with L_1 , L_3 must overlap with L_4 , L_1 must overlap with L_3 , and L_2 must overlap with L_3 .

3.2.2 Constructing and Coloring the Lineage Interference Graph

Next, we construct a lineage interference graph (LIG), an undirected graph whose vertices represent the lineages. Two vertices are connected by an interference edge if and only if the live ranges of the lineages represented by them definitely overlap. Using Theorem 3.1, the lineage interference graph can be computed from the transitive closure of the augmented DDG. Fortunately, the reach relation defined earlier significantly reduces the cost of generating the LIG by requiring only the verification of the paths between the nodes that start and the nodes that end lineages.

There is an edge in the LIG between two lineages $L_u = [u_1, u_2, \dots, u_m]$ and $L_v = [v_1, v_2, \dots, v_n]$ if and only if $R(u_1, v_n) = 1$ and $R(v_1, u_m) = 1$. In our motivating example, the set of start nodes is $S = \{a, c, d, e\}$ and the set of end nodes is $E = \{f, g, h\}$. For the presentation of our algorithm, the reach relation can be represented by a binary matrix where each row is associated with a node in S and each column is associated with a node in E . The reach relation for this example is shown in Fig. 4a. The LIG for this example is shown in Fig. 2c.

The lineage interference graph can be colored using a heuristic graph coloring algorithm [10], [11]. We refer to the number of colors required to color the interference graph as the Heuristic Register Bound (HRB). It should be noted that, due to the heuristics involved in coloring the interference graph and due to the sequencing order of descendant nodes

in the DDG, the HRB computed is a near-optimal solution. For the motivating example shown in Fig. 1a, the lineage interference graph is depicted in Fig. 2c. This interference graph can be colored using three colors, i.e., the HRB for our motivating example is 3.

3.3 Lineage Fusion

Before we proceed to describe our instruction sequence generation method, in this section, we present an optimization that *fuses* two lineages into a single one. Lineage fusion increases the success of the instruction sequence generation algorithm on finding an instruction sequence that requires no more than HRB registers.

3.3.1 Conditions for Lineage Fusion

Let $L_u = [u_1, u_2, \dots, u_m]$ and $L_v = [v_1, v_2, \dots, v_n]$. If u_1 reaches v_n , but v_1 does not reach u_m , then the live ranges of L_u and L_v do not necessarily interfere with each other, i.e., it is possible to generate a legal sequence where these two live ranges do not overlap. In this case, it is possible to use the same register for both lineages. However, if we start listing the nodes of L_v in the instruction sequence before we finish listing all the nodes of L_u , the two lineages will interfere and we will not be able to use the same register for both lineages. An instruction sequencing method that uses a fixed number of (say, HRB) registers will deadlock in such a situation because it will not be able to release the register assigned for the nodes of L_u for usage by L_v . This lineage interference, and consequent sequencing deadlock, is avoidable. Hence, we call this deadlock an *avoidable deadlock*. Further discussion on deadlocks is deferred to Section 3.4.

A simple solution that prevents the above interference and the consequent sequencing deadlock involves introducing a new sequencing constraint in the DDG that forces all the nodes of L_u to be listed before any node of L_v is listed. We call this operation a *lineage fusion* because it treats lineages L_u and L_v as if they were a single lineage. An additional advantage of lineage fusion is a reduction in the number of nodes in the lineage interference graph and a consequent reduction in the cost of coloring that graph. After the lineages are formed as described in Section 3.1, pairs of lineages that satisfy the *Lineage Fusion Condition* can be fused into a single lineage. Formally,

Definition 3.6. Two lineages $L_u = [u_1, u_2, \dots, u_m]$ and $L_v = [v_1, v_2, \dots, v_n]$ can be fused into a single lineage if:

1. u_1 reaches v_n , i.e., $R(u_1, v_n) = 1$;
2. v_1 does not reach u_m , i.e., $R(v_1, u_m) = 0$.

When the lineages L_u and L_v are fused together, a sequencing edge from u_m to v_1 is inserted in the DDG. The lineages L_u and L_v are removed and a new lineage $L_w = [u_1, u_2, \dots, u_m] \cup [v_1, v_2, \dots, v_n]$ is inserted in the lineage set. Note that the last node of the first lineage, u_m , does not necessarily use the same registers as the other nodes in the new L_w lineage. Thus, it is important to represent the lineage resulting from the fusion as a union of semi-open sequences as we did in this paragraph. Fusing two lineages L_u and L_v causes the respective nodes u and v in the interference graph to be combined into a single node, say w . Every edge incident on u or v is now incident on w . In

fusing lineages, the reach relation should be updated after each fusion. We discuss how to update the reach relation in the following subsection. Before that we establish that the sequencing edges introduced by lineage fusion do not result in cycles in the DDG.

Lemma 3.2. *The fusion of lineages does not introduce any cycle in the augmented DDG.*

Proof. We use proof by contradiction. Let $L_u = [u_1, u_2, \dots, u_m]$ and $L_v = [v_1, v_2, \dots, v_n]$ be two lineages such that 1) u_1 reaches v_n and 2) v_1 does not reach u_m . Let us assume that, when we fuse L_u and L_v , the inclusion of the sequencing edge from u_m to v_1 causes the formation of a cycle in the DDG. If this edge (u_m, v_1) forms a cycle, then this cycle must include a path from v_1 to u_m . But, that contradicts condition 2 for fusion. Hence, the fusion of two lineages under the conditions of Definition 3.6 cannot introduce a cycle in the augmented DDG. \square

Lineage fusion helps to reduce the number of partially overlapping live range pairs and thereby reduces the register requirement. It accomplishes this by fusing the two lineages corresponding to the partially overlapping live ranges into one and forcing an instruction sequence order on them. Lineage fusion is applied after lineage formation and before the coloring of the lineage graph. Therefore, the interference graph to be colored after lineage fusion has fewer vertices. It would be also legal to fuse the two lineages L_u and L_v when $R(v_1, u_m) = R(u_1, v_n) = 0$. However, such fusions would impose an unnecessary constraint in the sequencing order of L_u and L_v restricting the freedom of the sequencing algorithm. The sharing of registers by completely independent lineages is indicated by the coloring of the LIG.

3.3.2 Implementing Lineage Fusion

While it may seem at first that it would be necessary to recompute the transitive closure of the augmented DDG after each lineage fusion, we can obtain all possible lineage fusions from the reach relation already computed for the construction of the lineage interference graph.

To find candidates for fusion, we search for a pair of lineages $L_u = [u_1, u_2, \dots, u_m]$ and $L_v = [v_1, v_2, \dots, v_n]$ such that $R(u_1, v_n) = 1$ and $R(v_1, u_m) = 0$. The fusion condition says that we can create the new lineage $L_w = L_u \cup L_v$. The creation of this new lineage will require the addition of a scheduling edge from u_m to v_1 . The algorithm updates the reach relation in such a way that all the nodes that could reach u_m before the fusion, can also reach v_n now. In other words, $R(u_1, v_b) = 1$, for some $v_b \in E$, if either $R(u_1, v_b) = 1$ before the lineage fusion, or $R(v_1, v_b) = 1$. After the fusion, v_1 can no longer be in the set S of nodes that start lineages. Notice that v_1 is no longer the start node of the fused lineage L_w and that each node can start at most one lineage in the lineage formation algorithm. Therefore, the row corresponding to v_1 should be eliminated from the reach relation. Further, if u_m does not terminate any other lineage, its

column can also be eliminated from the reach relation.⁶ Furthermore, node v_1 is removed from S . Likewise, node u_m can be removed from E , the set of nodes that end a lineage, if u_m does not end any lineage other than L_u .

When there are multiple candidates for lineage fusion, we arbitrarily select a pair of lineages to be fused. Fusion results in updating the sets S and E and the reach relation. We keep searching the reach relation representation for new pairs of lineages that can be fused until none is found. For instance, a lineage L_u may be first fused with a lineage L_v to create a new lineage $L_u \cup L_v$ and then this compound lineage may be fused with a third lineage L_w to form $L_u \cup L_v \cup L_w$. We note here that the order in which lineages are fused may result in different HRB values. A study on the order of lineage fusion and its impact is beyond the scope of this paper.

In our motivating example, the initial analysis of the reach relation finds that lineage $L_4 = [d, g]$ can be merged with lineage $L_2 = [c, f]$ because $R(d, f) = 1$ and $R(c, g) = 0$. The new lineage is $L_4 - 2 = [d, g] \cup [c, f]$ and the new sequencing edge in the DDG is (g, c) . The updated reach relation is shown in Fig. 4b. Because there are no more zeros in the reach relation, no more lineage fusions are possible. The DDG and the LIG after fusion are shown in Fig. 4c and Fig. 4d.

In the following subsection, we present a heuristic method to derive an instruction sequence for the augmented DDG.

3.4 Instruction Sequence Generation

The coloring of the lineage interference graph associates a register with each lineage and, hence, with the nodes in the lineage. For example, a register assignment A with three colors for the nodes in our motivating example is:

$$A = \{(a, R_1); (b, R_1); (f, R_1); (c, R_2); (d, R_2); (e, R_3); (g, R_3)\},$$

where R_1 , R_2 , and R_3 are general purpose registers. Our solution assumes that registers that are live-in and live-out in the DDG will be assigned by a global register allocation procedure. However, as discussed in Section 2, our method can account for live-in and live-out registers through the addition of dummy *source* and *sink* nodes.

Our sequencing method is a modified list scheduling algorithm. It uses the register allocation information obtained by coloring the lineage interference graph. The sequencing method *lists* nodes from a ready list based on the height priority and on the availability of registers assigned to them.

The instruction sequence generation algorithm is shown in Fig. 5. The sequencing algorithm takes as inputs G' , the DDG augmented with sequencing edges; L , a list of lineages obtained from the lineage formation algorithm with lineage fusion applied, and A , the register assignment for the nodes from the coloring of the lineage interference graph. The availability of the required register needs to be checked only if node v_i is the first node in its lineage. Otherwise, the predecessor of v_i will pass the register assigned to it, which

```

SEQUENCING( $G', L, A$ )
1. ReadyList  $\leftarrow \{(v_i, R_j) \text{ such that } v_i \text{ has no predecessors}\}$ 
2. RegAvailable  $\leftarrow \{R_1, R_2, \dots, R_N\}$ 
3. while ReadyList  $\neq \emptyset$  do
4.   for each node  $v_i$  in the ReadyList in decreasing height order do
5.     if ( $v_i \notin S$ ) or ( $((v_i, R_j) \in A)$  and ( $R_j \in \text{RegAvailable}$ ))
        // either  $v_i$  is not the start node of a lineage or
        // the register assigned to  $v_i$  in  $A$  is available
6.       Remove  $R_j$  from RegAvailable
7.       Remove  $(v_i, R_j)$  from ReadyList
8.       List( $v_i, R_j$ )
9.       Add to the ReadyList all successors of  $v_i$ 
        that have all its predecessors listed
10.      if ( $v_i \in E$ ) and  $(v_i, R_j) \in A$ 
        // node  $v_i$  ends a lineage which
        // is assigned register  $R_j$ 
11.        Return  $R_j$  to RegAvailable
12.      endif
13.    endif
14.  endfor
15. end while

```

Fig. 5. Sequencing Algorithm.

will be free when v_i is ready to be listed. The algorithm uses the augmented DDG with sequencing edges, which ensure that v_i is the *last use* of its predecessor. When two lineages $L_u = [u_1, u_2, \dots, u_m]$ and $L_v = [v_1, v_2, \dots, v_n]$ are fused together, they form a single lineage L_w and a single register is assigned to the nodes $u_1, u_2, \dots, u_{m-1}, v_1, \dots, v_{n-1}$ of lineage L_w . As lineage fusion introduces a sequencing edge from u_m to v_1 and removes nodes u_m and v_1 from E and S , respectively, there is no need to check for the availability of a register when v_1 is listed.

Unfortunately, the above sequencing algorithm may result in a deadlock due to two reasons [19]. First, the order of listing two nodes belonging to two different lineages that are assigned the same color may result in a deadlock. We refer to the deadlocks caused by the wrong ordering of nodes as *avoidable deadlocks*, as they can often be avoided through the lineage fusion process. We illustrate this with the help of an example in the following discussion. The second kind of deadlocks, referred to as *unavoidable deadlocks*, are caused due to the underestimation of HRB. This could happen because the condition used to test if two live ranges definitely overlap (stated in Theorem 3.1) is sufficient but not necessary. As illustrated below, lineage fusion helps to reduce the occurrences of both avoidable and unavoidable deadlocks.

Applying the sequencing algorithm to our motivating example, we first list node a . If lineage fusion is not used, the listing of node a causes nodes c , d , and e to be added to the Ready List. Since register R_3 is available, node e can be listed next. The sequencing algorithm is left with only register R_2 . Now, either node d or node c can be listed. As there is no criterion to choose between them, the tie is broken arbitrarily. Unfortunately, if node c is listed next, a cycle of dependences is created because node d cannot be listed until node f is listed, which would release register R_2 . However, in order to list f , we must list b before that. But, b cannot be listed before d as b must be the last use of R_1 . This creates a cycle and the sequencing algorithm deadlocks. This is an avoidable deadlock that can be solved by lineage fusion.

6. A node might end multiple lineages because the last node of the lineage does not share a register with other nodes in the lineage. For instance, in our motivating example node, h ends lineages L_1 and L_3 .

If lineage fusion is employed, lineages $L4$ and $L2$ are fused together and a sequencing edge from node g to node c is added to the graph. This sequencing edge (g, c) ensures that node d is listed before c . Thus, nodes d, g, c, b, f , and h are listed subsequently in that order. The instruction sequence requires only three registers, as shown in Fig. 1c.

Unfortunately, even with the application of lineage fusion, *unavoidable deadlocks* occur when the heuristic register bound (HRB) computed from coloring the lineage interference graph is lower than the actual number needed. In this case, there does not exist a legal instruction sequence that uses HRB or fewer registers. To overcome the deadlock problem (both avoidable and unavoidable deadlocks), our sequencing algorithm follows a simple heuristic that increases the HRB by one. The algorithm then picks one of the nodes (the one with the maximum height) in the Ready List and changes its register assignment (as well as that of the remaining nodes in that lineage) to a new register. This strategy overcomes a deadlock by gradually increasing the HRB and trying to obtain a sequence that uses as few extra registers as possible. We measure the performance of our heuristic approach in Section 5.

4 EXACT APPROACH TO THE MRIS PROBLEM

In this section, we formulate the MRIS problem as an integer linear programming (ILP) problem. Our ILP formulation obtains, for each node u in the DDG, a unique assignment f_u that represents u 's position in the instruction sequence. If there are n nodes in the DDG, then f_u can assume values from 1 to n . That is,

$$0 < f_u \leq n. \quad (1)$$

Further, since we are interested only in a sequential ordering of the nodes, no two nodes can have the same position, i.e., $f_u \neq f_v$ for all $u \neq v$. We represent this relation as an integer linear constraint using the technique proposed by Hu [22]. We use a 0-1 integer variable $w_{u,v}$ and write the condition $f_u \neq f_v$ as a pair of constraints:

$$f_u - f_v > -n * w_{u,v}, \quad (2)$$

$$f_v - f_u > -n * (1 - w_{u,v}), \quad (3)$$

where n is the number of nodes in the DDG. Intuitively, $w_{u,v}$ represents the sign of $(f_v - f_u)$ and $w_{u,v}$ is 1 if the sign is positive and 0 otherwise.

Next, we derive conditions that ensure that the instruction sequence obeys true data dependences. If there exists a flow dependence arc (u, v) in the DDG, then clearly node u must be listed ahead of node v . That is,

$$f_v - f_u > 0 \quad \text{for all } (u, v) \text{ in the DDG.} \quad (4)$$

Assume that node u has k successors: v_1, v_2, \dots, v_k . The live range of the value produced by node u is the interval $[f_u + 1, \max(f_{v_1}, f_{v_2}, \dots, f_{v_k})]$. Thus, the value produced by node u is live at position t if and only if 1) $f_u + 1 \leq t$ and 2) $(f_{v_1} \geq t) \vee (f_{v_2} \geq t) \vee \dots \vee (f_{v_k} \geq t)$. For t ranging from 1 to n , the value of the 0-1 integer variable $d_{u,t}$ represents whether $f_u \leq t$ is true (see inequalities (5) and (6)). Similarly, the value of the 0-1 integer variable $e_{i,u,t}$ indicates

whether $f_{v_i} > t$ (see inequalities (7) and (8)). The constraints on the position of all nodes that have true dependences on node u , i.e., nodes v_1, \dots, v_k are represented in a similar way.

$$t - f_u \geq -n * (1 - d_{u,t}) \quad (5)$$

$$t - f_u < n * d_{u,t} \quad (6)$$

$$t - f_{v_i} \geq -n * e_{i,u,t} \quad (7)$$

$$t - f_{v_i} < n * (1 - e_{i,u,t}). \quad (8)$$

The 0-1 integer variable $e_{u,t}$ is 1 if and only if any of the $e_{1,u,t}, \dots, e_{k,u,t}$ variables is 1. This condition is imposed by the following two constraints:

$$e_{1,u,t} + \dots + e_{k,u,t} \geq e_{u,t} \quad (9)$$

$$e_{1,u,t} + \dots + e_{k,u,t} \leq k * e_{u,t}. \quad (10)$$

Finally, the 0-1 integer variable $s_{u,t}$ indicates whether variable u is live at time t . Thus, $s_{u,t}$ is 1 if and only if both $d_{u,t}$ and $e_{u,t}$ are nonzero. This relation is represented by the following two constraints:

$$d_{u,t} + e_{u,t} \geq 2 * s_{u,t} \quad (11)$$

$$d_{u,t} + e_{u,t} \leq 2 * s_{u,t} + 1. \quad (12)$$

The interference relations between the live ranges within a basic block form an acyclic interval graph. Therefore, it is optimally colorable with the same number of colors as the maximum width of the interval graph [21]. The width of the interval graph at position t can be represented by $s_{1,t} + s_{2,t} + \dots + s_{n,t}$. Hence, the objective function is to **minimize** z , where

$$z \geq \sum_{i=1}^n s_{i,t} \quad \text{for all } t \in [1, n]. \quad (13)$$

Thus, the ILP problem is to minimize z subject to inequalities (1) to (13).

5 EXPERIMENTAL RESULTS

In this section, we present results from two sets of experiments designed to evaluate the performance of our lineage based algorithm. In the first set of experiments (see Section 5.1), we compare the register requirement of the instruction sequence generated by our lineage-based sequencing algorithm with the minimum register requirement computed by the ILP formulation presented in Section 4. Because of the time complexity of the ILP algorithm, the ILP method can find a solution in a reasonable time only for DDGs with a limited number of nodes (in spite of the use of an efficient commercial ILP solver, viz CPLEX). Therefore, instead of SPEC benchmarks, we use a set of loops extracted from a collection of benchmark programs for this comparison. We were able to compare the results for 675 DDGs and determined that, in 99.2 percent of them, our sequencing

TABLE 1
Comparison between the Register Requirement of ILP (Exasct) Approach, HRB (Estimation),
and a Heuristic MRIS Sequence for MOST DDGs

Estimate's Relation	Register Requirement			Number of DDGs		
	ILP	HRB	Sequence	DDGs	Total	Percentage
ILP = HRB = Sequence	1-8	1-8	1-8	650	650	96.3%
ILP = Sequence > HRB	3	2	3	10	19	2.8%
	4	3	4	5		
	5	4	5	1		
	7	3	7	1		
	8	7	8	2		
ILP = HRB < Sequence	4	4	5	2	3	0.4%
	6	6	7	1		
ILP < HRB = Sequence	2	3	3	1	3	0.4%
	3	4	4	2		

method found a sequence that uses the minimum number of registers.

For our second set of experiments, we implemented our method in the SGI MIPSpro compiler suite, a set of highly optimizing compilers for Fortran, C, and C++ on MIPS processors. We report both static and dynamic performance measures for four different versions of the compiler on the SPEC95 floating-point suite.⁷ The static performance measures include the number of basic blocks that required register spilling, the average spills per basic block as well as the total number of spill instructions inserted in the code. We report execution time of the compiled program on a MIPS R10000 and the number of dynamic loads and stores that graduate as dynamic performance measures. The static and dynamic performance measures are reported in Section 5.2.

5.1 Comparing the Heuristic with the Exact Results

In order to compare the heuristic register bound (HRB) found by our method with the exact register requirement obtained with the ILP approach, we use a set of loops from the MOST framework [2]. The MOST framework consists of a collection of 1,200 loops extracted from SPEC92 (integer and floating-point), linpack, livermore, and the NAS kernels benchmarks. These DDGs were extracted using an optimizing research compiler. However, many of these DDGs had a large number of nodes and edges; when the exact solution was tried for these DDGs, the ILP solver couldn't get a solution, even after a reasonably long execution time. In the interest of time, we eliminated such DDGs and used the remaining 675 DDGs in our experiments. The DDGs considered in our experiments vary widely in size with a median of 10 nodes, a geometric mean of 12 nodes, and an arithmetic mean of 19 nodes per DDG.

Out of the 675 DDGs, in 650 (96.3 percent of the total) the estimated HRB is exact and our sequencing algorithm obtains an instruction sequence that uses the minimum number of registers. In Table 1, we present the exact

minimum register requirement found by the ILP formulation, the estimated register requirement found by the HRB approach, and the number of registers used by our lineage-based sequencing algorithm. We separate the loops in which these estimates are identical to the minimum register requirement (the 650 loops in the first row of the table) from those in which they are not (the last three rows in the table). In 19 loops (2.8 percent of total), although the HRB is an underestimation of the number of registers required, the sequencing algorithm based on lineage coloring used the minimum number of registers required. In three DDGs (0.4 percent of total), although the HRB is a correct estimation of the minimum number of registers, the sequencing algorithm used one extra register. And, in the other three loops, both the HRB estimate and the number of registers required by our heuristic sequencing algorithm are one more than the optimal value. Thus, out of the 675 DDG tested, only in six of them (0.8 percent) does the heuristic sequencing use one more register than the minimum required. These results give us confidence in the quality of the instruction sequences produced by our heuristic approach.

5.2 Comparison with a Production Compiler

In the second set of experiments, we implemented our heuristic instruction sequencing method in the SGI MIPSpro compiler. This compiler performs extensive optimizations including copy propagation, dead-code elimination, if-conversion, loop unrolling, cross-iteration optimization, recurrence breaking, instruction scheduling, and register allocation. The compiler also implements an integrated global-local scheduling algorithm [26] that is performed before and after register allocation. The global register allocation, based on graph coloring [8], [12], is followed by local register allocation. After register allocation, the data dependence graph is rebuilt and a postpass scheduling is invoked. During local register allocation, the MIPSpro compiler inserts spill code according to a cost function that estimates the number of spill load/store instructions that each spill candidate will incur. This estimation takes into consideration only the uses within the same basic block.

In our implementation, the algorithm described in Section 3 is used to optimize the instruction sequence at the basic block level. This local optimization is applied only to basic

7. Our experiments with the SPEC Int suite resulted in no measurable differences in execution time between the HRB and the Optimized versions of the compiler. This is possibly due to the facts that the basic blocks in the Spec Int benchmarks are smaller in size and our implementation of the MRIS approach performs instruction sequencing only within the basic block.

blocks that require spill code under the initial local register allocation. After the instruction sequence is optimized, the local register allocation is invoked again on the new instruction sequence. The goal of the instruction sequencing algorithm is to reduce the amount of spill code executed by reducing the register pressure of the generated code. We refer to this version of the compiler as HRB-based sequencing, or simply as the HRB approach.⁸

The performance of the HRB version is evaluated against a baseline version of the MIPSpro compiler. We also measure the HRB approach against an optimized version of the MIPSpro compiler which includes a combined instruction scheduling and register allocation algorithm. Thus, in the experimental results presented in this section, we compare four versions of the compiler:

Baseline: The Baseline compiler includes several traditional optimizations, such as copy propagation, dead-code elimination, etc., listed earlier in this subsection. In fact, all four compiler versions apply these optimizations. In addition, the Baseline compiler traverses the instructions of a basic block in reverse order to perform local register allocation, but does not try to optimize the instruction sequencing when the local register allocator requires spill code.

Optimized: In this version of the compiler, in addition to the traditional optimizations, there is a combined instruction scheduling and register allocation algorithm that is implemented in the MIPSpro compiler. The instruction sequencing algorithm used for this optimization is a depth-first traversal algorithm that takes resource constraints into consideration.

HRB: The HRB compiler version also includes traditional optimizations, but, when the local register allocator detects the need to introduce spill code in a basic block, our instruction sequencing algorithm based on the formation of lineages (including lineage fusion) and on the heuristic register bound described in this paper is applied. Register allocation is performed in the new instruction sequence.

HRB (No Fusion): Identical to HRB, except that the lineage fusion algorithm is not implemented.

The Baseline compiler is an optimized implementation including the integrated global and local scheduling, global and local register allocation, and optimized spilling. However, when register pressure is high (which is the case for some of the SPEC FP benchmarks), a more sophisticated instruction scheduling geared toward reducing register pressure is needed. The Optimized version of the compiler uses this additional instruction scheduling, while the Baseline does not.

We present our performance results for a machine with 32 integer and 32 floating-point architected registers.

8. It should be noted that, although the original objective of the MRIS problem is to arrive at an instruction sequence that uses the minimum number of registers for a basic block, in our implementation, the HRB-based sequencing is applied only to those basic blocks whose register requirement is greater than the available number of registers. Thus, the HRB sequencing approach is only used to reduce the register pressure for basic blocks which do incur register spills.

Because the HRB algorithm is more effective in applications with high register pressure, we will also present the performance results for a machine with 32 integer and 16 floating-point architected registers. This is the same target processor, but, in the machine description file in the compiler, we restrict the available FP registers to 16 so that the compiler cannot use half of the FP register file. This in effect increases the register pressure and, hence, spills. These results are indicative of the performance of our method on applications with higher register pressure.

5.2.1 Static Measurements

First, we report static performance measures such as the number of spill instructions inserted in the code by the different versions of the compiler. Columns 2, 4, 7, and 10 in Table 2 report the number of basic blocks in which spill code was introduced by the different versions of the compilers.⁹ We report the total number of spill instructions introduced by the different versions of the compiler in each application in columns 3, 5, 8, and 11. Last, columns 6, 9, and 12 represent, respectively, the percentage reduction in the number of (static) spill instructions introduced by the Optimized, HRB (No Fusion), and HRB versions of the compiler compared to the Baseline version. We define the percentage reduction in the spill instructions compared to the Baseline version of the compiler as:

$$\% \text{ Reduction}(P, V) = 100 \times \frac{(S(P, B) - S(P, V))}{S(P, V)}, \quad (14)$$

where $S(P, V)$ is the number of spills in benchmark P under version V , V can be the Optimized, HRB without fusion, or HRB version of the compiler, and $S(P, B)$ is the number of spills in benchmark P under the Baseline version. The top or bottom section of the table reports the performance for a target machine with 32 or 16 FP registers. In each section, the last row reports the average number of basic blocks that incurred spill, the average number of spill instructions per benchmark, and the percentage reduction in the average spill instructions compared to the Baseline version. The average percentage reductions in this row are computed as the percentage reduction in the average number of spills inserted in the benchmarks and not as the average of the percentage reduction for the different benchmarks.

In Table 2, the total number of basic blocks that require spilling for a specific version of the compiler can be obtained by adding the entries in the column "Blocks with Spill." For a machine with 16 FP registers, the total number of blocks that require any spilling is reduced from 210 in the Baseline compiler to 94 in the HRB compiler, i.e., 55 percent of the blocks that required spill operations before no longer do. Compared to this, the Optimized MIPSpro version inserts spill code in 152 blocks (a reduction of only 28 percent). For a target machine with 32 FP registers, the Baseline compiler inserts spills in 88 blocks compared to 32 by the HRB version

9. Besides the SPEC95 floating-point benchmarks reported in the paper, we also conducted our experiments for *swim*, *mgrid*, and *hydro2d*, but, for these benchmarks, there are no basic blocks with spills, even under the Baseline version of the compiler. Hence, we do not include them in our results.

TABLE 2
Static Count of the Number of Spill Operations Inserted

Benchmark	Compiler Version										
	Baseline		Optimized			HRB (No Fusion)			HRB		
	Blocks w/spills	Total Spills	Blocks w/spills	Total Spills	% reduc.	Blocks w/spills	Total Spills	% reduc.	Blocks w/spills	Total Spills	% reduc.
A Machine with 32 INT and 32 FP Registers											
tomcatv	1	1	0	0	100.00	0	0	100.00	0	0	100.00
su2cor	8	36	1	6	83.33	1	5	86.11	1	5	86.11
applu	22	1080	14	454	57.96	11	337	68.80	10	340	68.52
turb3d	10	67	7	69	-2.99	1	28	58.21	0	0	100.00
apsi	20	287	11	135	52.96	9	119	58.54	9	113	60.63
fpppp	17	1180	13	782	33.73	12	818	30.68	12	560	52.54
wave5	10	111	5	50	54.95	1	1	99.10	0	0	100.00
Average	12.57	394.57	7.29	213.71	45.84	5.00	186.86	52.64	4.57	145.43	63.14
A Machine with 32 INT and 16 FP Registers											
tomcatv	2	38	1	22	42.11	2	16	57.89	1	13	65.79
su2cor	13	80	2	22	72.50	5	44	45.00	1	5	93.75
applu	63	1678	51	1130	32.66	43	945	43.68	35	819	51.19
turb3d	23	530	18	342	35.47	14	325	38.68	10	169	68.11
apsi	42	694	30	332	52.16	25	213	69.31	17	146	78.96
fpppp	23	1939	18	1391	28.26	16	1540	20.58	14	1276	34.19
wave5	44	827	32	492	40.51	26	223	73.04	16	124	85.01
Average	30.00	826.57	21.71	533.00	35.52	18.71	472.29	42.86	13.43	364.57	55.89

(a reduction of 63.6 percent) and 51 blocks by the Optimized version (a reduction of 42.1 percent).

A comparison of the average number of spills in each benchmark, reveals that the HRB version reduces the number of spills by 63.1 percent and 55.9 percent, respectively, in the 32 FP and 16 FP register machines, compared to the Baseline compiler. More specifically, for a machine with 16 FP registers, in terms of the total number of spills, the reduction due to HRB ranges from 34 percent to 94 percent. For a machine with 32 FP registers, the percentage reduction varies from 52.5 percent to 100 percent. In comparison, the Optimized version reduces the average number of spills by only 45.8 percent and 35.5 percent, respectively for the two machines. The most dramatic improvement when comparing the HRB version with the Optimized compiler is observed for fpppp, applu, and wave5 benchmarks, for a machine with 32 and 16 FP registers, where the Optimized version incurs 200 or more static spill instructions compared to the HRB approach. Thus, even compared to the Optimized compiler, the HRB version produces code that significantly reduces the number of spills.

The fusion of lineages reduces the number of spill operations inserted in the code in relation to a version of the HRB algorithm that does not perform lineage fusion. Compared to the HRB version, which reduced the average spills by 63.1 percent and 55.9 percent, respectively, for machines with 32 and 16 FP registers, the HRB version without lineage fusion resulted in a reduction of only 52.6 percent and 42.9 percent.

Last, notice that the reduced percentage improvement in average spills for the 32 FP register machine is less than that for a machine with 16 FP registers. This result may seem to be counterintuitive, given that the HRB algorithm should

perform well for programs with higher register pressure. Note, however, that, in a machine with 16 FP registers, 462 spills were eliminated on an average, while, for a machine with 32 FP registers, this number was 249. Thus, the benefit due to the HRB approach is larger in benchmarks with higher register pressure. The lower percentage reduction observed is only due to the large denominator value (average spills in the Baseline version).

5.2.2 Dynamic Measurements

To report dynamic performance, we conducted our experiments on a Silicon Graphics machine with a 194 MHz MIPS R10000 processor, 32 KB instruction cache, 32 KB data cache, 1 MB of secondary unified instruction/data cache, and 1 GB of main memory. We measured the wall clock time for the execution of each benchmark under the IRIX 6.5 operating system with the machine running in a single user mode. As the emphasis of our work is on sequencing the instructions to reduce the register requirements and spill code, we used the R10000 hardware counters and the perfex tool to measure the number of loads and stores graduated in each benchmark under the different versions of the compiler. Since the Baseline and HRB versions of the compiler are identical except for the instruction reordering at the basic block level, the reduction in the number of loads/stores executed in each benchmark program corresponds to the number of spill loads/stores reduced by the HRB approach.

First, we report the dynamic measurements for spill instructions. The number of loads and stores graduated from the pipeline for each benchmark under the different versions of the compiler are shown in Tables 3 and 4. Each table shows the number of operations (in billions) as well as the percentage reductions in relation to the Baseline

TABLE 3
Number of Graduated Loads (in Billions) and Percentage Reduction in Relation to the Baseline Compiler

Benchmark	Compiler Version						
	Baseline	Optimized		HRB (No Fusion)		HRB	
	# Loads	# Loads	% reduc.	# Loads	% reduc.	# Loads	% reduc.
A Machine with 32 INT and 32 FP Registers							
tomcatv	6.51	6.37	2.1	6.07	6.7	6.19	5.0
su2cor	5.83	5.82	0.1	5.83	0.1	5.82	0.1
applu	8.07	7.56	6.3	7.51	7.0	7.51	7.0
turb3d	11.63	11.63	0.0	11.58	0.4	11.59	0.3
apsi	4.96	4.83	2.5	4.73	4.6	4.67	5.7
fpppp	28.40	24.65	13.2	23.13	18.6	22.45	20.9
wave5	3.72	3.72	-0.2	3.69	0.8	3.69	0.7
Average	9.87	9.23	6.53	8.93	9.53	8.85	10.4
A Machine with 32 INT and 16 FP Registers							
tomcatv	8.42	7.84	6.9	7.25	13.9	6.88	18.3
su2cor	6.29	5.92	5.9	6.48	-3.0	5.68	9.7
applu	8.77	8.42	3.9	8.19	6.5	8.12	7.4
turb3d	14.75	13.30	9.8	12.02	18.5	12.65	14.2
apsi	5.38	5.15	4.4	5.02	6.8	4.70	12.7
fpppp	35.04	28.00	20.2	30.52	12.9	26.81	23.5
wave5	3.88	3.75	3.3	3.69	5.0	3.73	3.9
Average	11.79	10.34	12.3	10.45	11.33	9.8	16.9

TABLE 4
Number of Graduated Stores (in Billions) and Percentage Reduction in Relation to the Baseline Compiler

Benchmark	Compiler Version						
	Baseline	Optimized		HRB (No Fusion)		HRB	
	# Stores	# Stores	% reduc.	# Stores	% reduc.	# Stores	% reduc.
A Machine with 32 INT and 32 FP Registers							
tomcatv	2.34	2.34	0.0	2.34	0.0	2.34	0.0
su2cor	2.84	2.84	0.0	2.84	0.0	2.84	0.0
applu	5.24	4.75	9.3	4.65	11.3	4.65	11.2
turb3d	13.33	13.34	0.0	13.24	0.6	13.24	0.6
apsi	3.17	3.17	0.0	3.17	0.0	3.17	0.0
fpppp	20.78	20.68	0.5	19.67	5.3	19.55	5.9
wave5	3.76	3.76	0.0	3.76	0.0	3.76	0.0
Average	7.35	7.27	1.1	7.1	3.5	7.08	3.7
A Machine with 32 INT and 16 FP Registers							
tomcatv	3.22	2.63	18.2	2.54	21.2	2.73	15.2
su2cor	2.86	2.86	0.0	2.86	0.0	2.86	0.0
applu	5.88	5.37	8.6	5.31	9.7	5.16	12.3
turb3d	14.03	13.81	1.5	13.86	1.2	13.82	1.4
apsi	3.52	3.39	3.7	3.29	6.6	3.19	9.4
fpppp	21.37	22.09	-3.4	20.60	3.6	21.26	0.5
wave5	3.90	3.88	0.7	3.85	1.3	3.83	1.8
Average	7.83	7.72	1.4	7.47	4.5	7.55	3.5

compiler. The percentage reduction is computed using an equation similar to (14). The average number of loads reported in the last row of each section of the tables is simply the arithmetic mean of the number of graduated loads in each of the benchmarks. However, the percentage reductions in these rows is the percentage reduction of the average number of loads. In other words, the average reported under the “% reduction” column is not the average of the percentage reductions. The same averaging method is

used for reporting the graduated stores (see Table 4) and the execution time (see Table 5).

The HRB version reduces the average number of loads and stores executed, respectively, by 10.4 percent and 3.7 percent for a machine with 32 FP registers. These numbers for a machine with 16 registers are 16.9 percent and 3.5 percent, respectively. Even when compared with the Optimized MIPSpro compiler, the reduction in the average number of loads and stores is significant (almost

TABLE 5
Execution Time (in Seconds) and Reductions in the Execution Time in Relation to the Baseline Compiler

Benchmark	Compiler Version						
	Baseline	Optimized		HRB (No Fusion)		HRB	
	Time	Time	% reduc.	Time	% reduc.	Time	% reduc.
A Machine with 32 INT and 32 FP Registers							
tomcatv	358.99	354.67	1.20	360.38	-0.39	354.17	1.34
su2cor	230.12	229.99	0.06	230.37	-0.11	229.76	0.16
applu	397.41	392.16	1.32	388.59	2.22	385.92	2.89
turb3d	393.94	394.06	-0.03	393.62	0.08	395.17	-0.31
apsi	248.22	249.79	-0.63	249.53	-0.53	248.31	-0.04
fpppp	391.82	361.39	7.77	365.46	6.73	336.11	14.22
wave5	198.84	199.28	-0.22	198.50	0.17	198.81	0.01
Average	317.05	311.62	1.71	312.35	1.48	306.89	3.20
A Machine with 32 INT and 16 FP Registers							
tomcatv	375.00	363.31	3.12	362.55	3.32	358.58	4.38
su2cor	232.22	229.29	1.26	232.42	-0.09	231.73	0.21
applu	402.41	404.22	-0.45	397.28	1.27	393.65	2.18
turb3d	415.63	411.09	1.09	409.44	1.49	404.41	2.70
apsi	253.25	253.18	0.03	252.48	0.31	249.81	1.36
fpppp	447.98	411.82	8.07	414.27	7.53	396.79	11.43
wave5	200.43	198.79	0.82	200.55	-0.06	200.10	0.16
Average	332.42	324.53	2.37	324.14	2.49	319.30	3.95

400 and 200 million instructions, respectively). Thus, reducing the number of static spill instructions results in a reduction in the number of dynamic spill instructions executed. The drastic improvement in a single benchmark is observed for fpppp, where 8.2 and 5.95 billion loads are eliminated (compared to the Baseline version), for machines with 16 and 32 FP registers, respectively.

Next, we report whether the reduction in the number of loads and stores executed corresponds to a reduction in the execution time of the benchmarks. Table 5 presents the execution time for each benchmark under the different versions of the compiler. This is the wall-clock time, measured in seconds, required to execute each benchmark. We also present the percentage reduction in the execution time. First, we notice that, by and large, there is a correlation between the improvement in execution time in Table 5 and the reduction in the number of loads and stores graduated. For example, the drastic improvements in the number of loads and stores seen in fpppp do translate into a significant improvement in execution time: 14.2 percent and 11.4 percent, respectively, for machines with 32 and 16 FP registers. Once again, although the reduction in execution time compared to the Baseline version is similar for different target machine configurations for which they were compiled, the lower percentage reduction in the 16 FP register case is due to the larger execution time of the Baseline version.

From Table 5, it can be seen that the code produced by the HRB version of the compiler does, in fact, reduce the execution time, although the percentage reduction is somewhat low. We remark that one should not be discouraged by the somewhat low improvement in execution time due to the HRB approach. The HRB approach performs instruction sequencing solely with the objective of reducing the register requirement. Our implementation of the HRB

algorithm does not take into account the resource constraints of the architecture, often sacrificing some instruction-level parallelism. The average execution time under the HRB version is comparable to or better than that under the Optimized version in each of the benchmarks. Our argument to support this approach is that modern superscalar processors with their out-of-order issue hardware should be able to uncover the instruction-level parallelism that is obscured by false register dependences. Compared to this, the MIPSpro Optimized version uses an integrated register allocation and instruction scheduler. Our experiments, in fact, demonstrate that the HRB sequencing approach minimizes the register requirement. At the same time, the execution of these programs resulted in a performance that is on a par with or better than that for the Baseline or the Optimized versions in terms of execution time. We have also witnessed a significant reduction (a few billion) in the number of loads and stores. This reduction is primarily due to the savings in spill code. As indicated by Cooper and Harvey [13], reducing the spill load/store operations results in less cache pollution and, hence, higher cache performance. Further reducing memory operations can significantly reduce the power dissipated. These are other advantages of the HRB approach in addition to achieving a reasonable performance improvement in execution time.

6 RELATED WORK

Instruction scheduling [16], [28] and register allocation [1], [8], [10], [11], [15], [28], [30], [32], [39] are important phases in a high-performance compiler. The ordering of these phases and its implications on the performance of the code generated have been studied extensively for in-order issue superscalar processors and Very Long Instruction Word

(VLIW) processors. In such processors, it is often necessary to expose enough instruction-level parallelism even at the expense of increasing the register pressure and, to some extent, the amount of spill code generated. Integrated techniques that try to minimize register spills while focusing on exposing parallelism were found to perform well [5], [6], [27], [29], [31]. All these approaches work on a *given* instruction sequence and attempt to improve register allocation and/or instruction scheduling. In contrast, our MRIS approach generates an instruction sequence from a DDG where the precise order of instructions is not yet fixed.

Modern out-of-order issue superscalar processors, which have the capability to perform instruction scheduling and register renaming at runtime, shift the focus of the instruction-level parallelism compilation techniques. Studies on out-of-order issue processors indicate that reducing the register pressure and, hence, the number of memory spill instructions executed is more crucial to the performance than exposing higher instruction level parallelism at compile time [36], [40].

The Minimum Register Instruction Sequence (MRIS) problem studied in this paper is different from the traditional *register allocation problem* [1], [8], [10], [11], [15], [28], [30]. Recently, there have also been a few proposals on register allocation based on integer linear programming [3], [18], [25]. The input to the MRIS problem is a partially ordered sequence, specified by a DDG, instead of a totally ordered sequence of instructions. Although the absence of a total order of instructions makes the MRIS problem harder, it also enables the generation of an instruction sequence that requires less registers. The MRIS problem is also quite different from the traditional *instruction scheduling problem* [1], [16], [17], [28], [41]. In the traditional instruction scheduling problem, the main objective is to minimize the execution time (length) of the schedule, taking into account the execution latencies of each operation (instruction) in the DDG and the availability of function unit resources. This is in contrast to the MRIS problem, where only the true dependence constraints are observed. The MRIS problem is also closely related to the *optimal code generation (OCG)* problem [1], [34], [33]. An important difference between traditional code generation methods and our MRIS problem is that the former emphasizes reducing the code length (or schedule length) for a fixed number of registers, while the latter minimizes the number of registers used.

The unified resource allocator (URSA) method deals with function unit and register allocation simultaneously [4]. The method uses a three-phase *measure-reduce-assign* approach, where resource requirements are measured and regions of excess requirements are identified in the first phase. The second phase reduces the requirements to what is available in the architecture and the final phase carries out resource assignment. More recently, Berson et al. have proposed the use of *register reuse dags* for reducing the register pressure [5]. A register reuse dag is similar to a lineage discussed in this paper. They have evaluated register spilling and register splitting methods for reducing the register requirements in the URSA method for individual loops, rather than the whole application. Our work, in

contrast, reports static and dynamic performance measures on the SPEC95 floating-point benchmark suite.

In the experiments presented in Section 5, our lineage formation method was applied only to basic blocks for which the local register allocation resulted in spills. Thus, the aggressive sequentialization in our algorithm does not prevent parallelism when an allocation is found that does not incur spills. In an independent work, Touati proposes an approach to perform register allocation and code scheduling in a single pass by generating a schedule that maximizes the number of values live at the same time [38]. His reasoning is that, when more values are simultaneously live, the scheduler will find more opportunities to explore parallelism. Thus, instead of creating long chains (lineages in our case), he attempts to create maximal antichains. When maximal antichains result in a register requirement that is greater than the number of available registers, Touati uses an approximate algorithm to find a quasi-optimal set of serialization arcs to reduce the number of registers used. In the few examples that we examined, his register saturation reduction algorithm and our lineage fusion algorithm created the same set of sequencing edges. It is interesting to note that these two algorithms use quite contrasting approaches (maximizing vs. minimizing the register requirements) to achieve the same objective of minimizing the register spills.

Last, the lineage formation and the heuristic list scheduling methods proposed in this paper are major improvements over, respectively, the chain formation and the modified list scheduling method discussed in [19]. The chain formation method allocates, at least conceptually, one register for each arc in the DDG, and the set of chains must cover all arcs. That is, it must include each def-use, not just def-last_use, in a chain. Hence, the instruction lineage approach more accurately models the register requirement. Second, the lineage formation overcomes an important weakness of instruction chains, namely allocating more than one register for a node. Further, a number of heuristics, including lineage fusion, have been incorporated into the sequencing method to make it more efficient and obtain near-optimal solutions.

7 CONCLUSIONS

In this paper, we address the problem of generating an instruction sequence for a computation that is optimal in terms of the number of registers used by the computation. This problem is motivated by requirements of modern out-of-order issue processors. In out-of-order issue processors, excessive register spills can potentially degrade the performance and, hence, must be avoided, even at the expense of reducing the parallelism exposed at compile time. Another motivation for the MRIS problem stems from the fact that register spills lead to memory accesses, which are expensive in terms of power dissipation.

We proposed two solutions to the MRIS problem. First was a heuristic solution that uses lineage formation, lineage interference graph, and a modified and efficient list scheduling method to obtain an efficient near-optimal solution to the MRIS problem. The second approach is based on an elegant integer linear programming (ILP)

formulation. Compared to the exact ILP-based approach for MRIS, the heuristic approach results in optimal solution for 99.2 percent of the DDGs used in our experiments, although these DDGs are small in size and have low register requirements.

We evaluated the performance of our heuristic method by implementing it in the MIPSpro production compiler and running SPEC95 floating-point benchmarks. Our experimental results demonstrate that our instruction sequencing method, which attempts to minimize the register requirements, reduces the average number of basic blocks that require spilling by 62.5 percent and the average spill operations in the (static) code by 63.1 percent. As a consequence, the HRB approach also reduces the average number of loads and stores executed by 10.4 percent and 3.7 percent, respectively, for a machine with 32 integer and 32 floating-point registers. This reduction in loads and stores also results in an improvement in the average execution time by 3.2 percent.

The MRIS approach to register allocation and instruction sequencing might allow a smaller cache and reduce the traffic of data between the processor and the memory structure, thus contributing to the reduction of power consumption. It might also reduce the need for spilling from the register stack in an IA-64 machine. An experimental study to test these hypotheses is a possible future direction for research.

ACKNOWLEDGMENTS

This research is supported by the US National Science Foundation (NSF), by the US National Security Agency (NSA), by the US National Aeronautics and Space Administration (NASA), by the US Defense Advanced Research Projects Agency (DARPA), and by the Natural Sciences and Engineering Research Council (NSERC) of Canada. The authors would also like to acknowledge current and former members of CAPSL for valuable discussions and Jim Dehnert and Sun Chan, formerly with SGI, for discussions about the organization of the MIPSpro compiler. Special thanks to the anonymous reviewers whose comments helped them improve the quality of this manuscript.

REFERENCES

- [1] A.V. Aho, R. Sethi, and J.D. Ullman, *Compilers—Principles, Techniques, and Tools*, corrected ed. Reading, Mass.: Addison-Wesley, 1988.
- [2] E.R. Altman, "Optimal Software Pipelining with Function Unit and Register Constraints," PhD thesis, McGill Univ., Montréal, Québec, Oct. 1995.
- [3] A.W. Appel and L. George, "Optimal Spilling for CISC Machines with Few Registers," *Proc. ACM SIGPLAN 2001 Conf. Programming Language Design and Implementation*, pp. 243-253, June 2001.
- [4] D. Berson, R. Gupta, and M.L. Soffa, "URSA: A Unified ReSource Allocator for Registers and Functional Units in VLIW Architectures," *Proc. Conf. Parallel Architectures and Compilation Techniques (PACT '98)*, June 1998.
- [5] D. Berson, R. Gupta, and M.L. Soffa, "Integrated Instruction Scheduling and Register Allocation Techniques," *Proc. 11th Int'l Workshop Languages and Compilers for Parallel Computing*, Aug. 1998.
- [6] D.G. Bradley, S.J. Eggers, and R.R. Henry, "Integrating Register Allocation and Instruction Scheduling for RISCs," *Proc. Fourth Int'l Conf. Architectural Support for Programming Languages and Operating Systems*, pp. 122-131, Apr. 1991.
- [7] P. Briggs, K.D. Cooper, and L. Torczon, "Rematerialization," *Proc. ACM SIGPLAN '92 Conf. Programming Language Design and Implementation*, pp. 311-321, June 1992.
- [8] P. Briggs, K.D. Cooper, and L. Torczon, "Improvements to Graph Coloring Register Allocation," *ACM Trans. Programming Languages and Systems*, vol. 16, no. 3, pp. 428-455, May 1994.
- [9] J.L. Bruno and R. Sethi, "Code Generation for a One-Register Machine," *J. ACM*, vol. 23, no. 3, pp. 502-510, July 1976.
- [10] G.J. Chaitin, "Register Allocation and Spilling via Graph Coloring," *Proc. SIGPLAN '82 Symp. Compiler Construction*, pp. 98-105, June 1982.
- [11] G.J. Chaitin, M.A. Auslander, A.K. Chandra, J. Cocke, M.E. Hopkins, and P.W. Markstein, "Register Allocation via Coloring," *Computer Languages*, vol. 6, pp. 47-57, Jan. 1981.
- [12] F.C. Chow and J.L. Hennessy, "The Priority-Based Coloring Approach to Register Allocation," *ACM Trans. Programming Languages and Systems*, vol. 12, no. 4, pp. 501-536, Oct. 1990.
- [13] K.D. Cooper and T.J. Harvey, "Compiler-Controlled Memory," *Proc. Architectural Support for Programming Languages and Operating Systems*, pp. 2-11, Oct. 1998.
- [14] G.R. Gao, L. Bic, and J.-L. Gaudiot, *Advanced Topics in Dataflow Computing and Multithreading*. IEEE CS Press, 1995.
- [15] L. George and A.W. Appel, "Iterated Register Coalescing," *Conf. Record 23rd ACM SIGPLAN-SIGACT Symp. Principles of Programming Languages*, pp. 208-218, Jan. 1996.
- [16] P.B. Gibbons and S.S. Muchnick, "Efficient Instruction Scheduling for a Pipelined Architecture," *Proc. SIGPLAN '86 Symp. Compiler Construction*, pp. 11-16, June 1986.
- [17] J.R. Goodman and W.-C. Hsu, "Code Scheduling and Register Allocation in Large Basic Blocks," *Conf. Proc. 1988 Int'l Conf. Supercomputing*, pp. 442-452, July 1988.
- [18] D.W. Goodwin and K.D. Wilken, "Optimal and Near-Optimal Global Register Allocation Using 0-1 Integer Programming," *Software—Practice and Experience*, vol. 26, no. 8, pp. 929-965, Aug. 1996.
- [19] R. Govindarajan, C. Zhang, and G.R. Gao, "Minimum Register Instruction Scheduling: A New Approach for Dynamic Instruction Issue Processors," *Proc. 12th Int'l Workshop Languages and Compilers for Parallel Computing*, Aug. 1999.
- [20] R. Govindarajan, H. Yang, J.N. Amaral, C. Zhang, and G.R. Gao, "Minimum Register Instruction Sequence Problem: Revisiting Optimal Code Generation for DAGs," *Proc. Int'l Parallel and Distributed Processing Symp.*, Apr. 2001.
- [21] L.J. Hendren, G.R. Gao, E.R. Altman, and C. Mukerji, "A Register Allocation Framework Based on Hierarchical Cyclic Interval Graphs," *The J. Programming Languages*, vol. 1, no. 3, pp. 155-185, 1993.
- [22] T.C. Hu, *Integer Programming and Network Flows*, p. 270. Addison-Wesley, 1969.
- [23] W.W. Hwu, S.A. Mahlke, W.Y. Chen, P.P. Chang, N.J. Warter, R.A. Bringmann, R.G. Ouellette, R.E. Hank, T. Kiyohara, G.E. Haab, J.G. Holm, and D.M. Lavery, "The Superblock: An Effective Technique for VLIW and Superscalar Compilation," *J. Supercomputing*, vol. 7, pp. 229-248, 1993.
- [24] Intel, *Intel IA-64 Architecture Software Developer's Manual*, Jan. 2000.
- [25] C.W. Kessler, "Scheduling Expression DAGs for Minimal Register Need," *Proc. Eighth Int'l Symp. Programming Languages: on Programming Languages: Implementations, Logics, and Programs (PLILP '96)*, pp. 228-242, Sept. 1996.
- [26] S. Mantripragada, S. Jain, and J. Dehnert, "A New Framework for Integrated Global Local Scheduling," *Proc. 1998 Int'l Conf. Parallel Architectures and Compilation Techniques*, pp. 167-174, Oct. 1998.
- [27] R. Motwani, K.V. Palem, V. Sarkar, and S. Reyan, "Combining Register Allocation and Instruction Scheduling," technical report, Courant Inst. of Math. Sciences, New York Univ., 1996.
- [28] S.S. Muchnick, *Advanced Compiler Design and Implementation*. San Francisco: Morgan Kaufmann, 1997.
- [29] B. Natarajan and M. Schlansker, "Spill-Free Parallel Scheduling of Basic Blocks," *Proc. 28th Ann. Int'l Symp. Microarchitecture*, pp. 119-124, Dec. 1995.
- [30] C. Norris and L.L. Pollock, "Register Allocation over the Program Dependence Graph," *Proc. ACM SIGPLAN '94 Conf. Programming Language Design and Implementation*, pp. 266-277, June 1994.
- [31] S.S. Pinter, "Register Allocation with Instruction Scheduling: A New Approach," *Proc. ACM SIGPLAN '93 Conf. Programming Language Design and Implementation*, pp. 248-257, June 1993.

- [32] M. Poletto and V. Sarkar, "Linear Scan Register Allocation," *ACM Trans. Programming Languages and Systems*, 1998.
- [33] T.A. Proebsting and C.N. Fischer, "Linear-Time, Optimal Code Scheduling for Delayed-Load Architectures," *Proc. ACM SIGPLAN '91 Conf. Programming Language Design and Implementation*, pp. 256-267, June 1991.
- [34] R. Sethi, "Complete Register Allocation Problems," *SIAM J. Computing*, vol. 4, no. 3, pp. 226-248, Sept. 1975.
- [35] R. Sethi and J.D. Ullman, "The Generation of Optimal Code for Arithmetic Expressions," *J. ACM*, vol. 17, no. 4, pp. 715-728, Oct. 1970.
- [36] R. Silvera, J. Wang, G.R. Gao, and R. Govindarajan, "A Register Pressure Sensitive Instruction Scheduler for Dynamic Issue Processors," *Proc. Conf. Parallel Architectures and Compilation Techniques (PACT '97)*, pp. 78-89, Nov. 1997.
- [37] J.E. Smith and G. Sohi, "The Microarchitecture of Superscalar Processors," *Proc. IEEE*, vol. 83, no. 12, pp. 1609-1624, Dec. 1995.
- [38] S.A.A. Touati, "Register Saturation in Superscalar and VLIW Codes," *Proc. 10th Int'l Conf. Compiler Construction*, pp. 213-228, Apr. 2001.
- [39] O. Traub, G. Holloway, and M.D. Smith, "Quality and Speed in Linear-Scan Register Allocation," *Proc. ACM SIGPLAN '98 Conf. Programming Language Design and Implementation*, pp. 142-151, June 1998.
- [40] M.G. Valluri and R. Govindarajan, "Evaluating Register Allocation and Instruction Scheduling Techniques in Out-of-Order Issue Processors," *Proc. Conf. Parallel Architectures and Compilation Techniques (PACT '99)*, Oct. 1999.
- [41] H.S. Warren Jr., "Instruction Scheduling for the IBM RISC System/6000 Processor," *IBM J. Research and Development*, vol. 34, no. 1, pp. 85-92, Jan. 1990.
- [42] W. Wulf, R.K. Johnson, C.B. Weinstock, S.O. Hobbs, and C.M. Geschke, *The Design of an Optimizing Compiler*. New York: American Elsevier, 1975.



R. Govindarajan received the BSc degree from Madras University in 1981 and the BE (electronics and communication) and PhD (computer science) degrees from the Indian Institute of Science in 1984 and 1989. He held postdoctoral research positions at the University of Western Ontario, London, Ontario, Canada and McGill University, Montreal, Québec, Canada. He was a faculty member in the Department of Computer Science, Memorial University of Newfoundland,

St. John's, Newfoundland, Canada, between 1994-1995. Since then, he has been with the Supercomputer Education and Research Centre and the Department of Computer Science and Automation, Indian Institute of Science, where he is now an associate professor. His current research interests are in the areas of compilation techniques for instruction-level parallelism, compilation techniques for DSP and embedded processors, distributed shared memory architectures, and cluster computing. He has more than 60 publications in these areas in international journals and refereed international conference proceedings. He is a senior member of the IEEE, and a member of the IEEE Computer Society, ACM, and ACM SIGARCH.



Hongbo Yang received the BS degree from Wuhan University, and the ME degree from the Institute of Computing Technology, Chinese Academy of Sciences (ICTCAS) in 1994 and 1997, respectively. He was a Practice Research Fellow at ICTCAS from 1997 to 1999. Currently, he is pursuing the PhD degree at the Department of Electrical and Computer Engineering, University of Delaware. His research interests include high performance and low power compilation techniques. He is a student member of the IEEE.



José Nelson Amaral received the PhD degree in electrical and computer engineering from the University of Texas at Austin in 1994 for his design and evaluation of a parallel computer architecture for symbolic computing. He received the ME degree from the Instituto Tecnológico de Aeronáutica, in São José dos Campos, SP, Brazil, in 1989, and the BE degree from Pontifícia Universidade Católica do Rio Grande do Sul (PUCRS), in Porto Alegre, RS, Brazil in 1987. He is an associate professor in the Department of Computing Science at the University of Alberta, Canada. After receiving his PhD, he returned to Brazil as a professor of electrical engineering at PUCRS. From January 1998 to June 2000, he was a postdoctoral associate with the Computer Architecture and Parallel Systems Laboratory (CAPSL) at the University of Delaware. His current research interests include compiler design and optimization, computer architecture, and high-performance computer systems. Professor Amaral is a member of the editorial board of the *IEEE Transactions on Computers* and serves as proceedings chair for the International Parallel and Distributed Processing Symposium and the International Conference on High Performance Computing. He is a senior member of the IEEE.



Chihong Zhang received the BS degree (computer science) in 1984 and the MS degree (computer science) in 1987 from Tsinghua University, Beijing, China. From 1987 to 1997, he was a faculty member in the Department of Computer Science and was an associate professor before he joined the CAPSL (Computer Architecture and Parallel Systems Laboratory) group as a visiting research scholar. He joined the compiler development team of Conexant Systems Inc. in 1999 and is now with Qualcomm Inc. His major research interests are now in compiler optimization techniques for DSP and embedded processors. He is a member of the IEEE.



Guang R. Gao received the SM and PhD degrees in electrical engineering and computer science from the Massachusetts Institute of Technology in 1982 and 1986, respectively. Currently, he is a professor in the Department of Electrical and Computer Engineering at the University of Delaware, where he has been the founder and leader of the Computer Architecture and Parallel Systems Lab. Prior to that, he was a faculty member of the School of Computer Science, McGill University, Montreal, Canada. His main research interests include high-performance computing systems and architectures, compilers, parallel programming, and applications such as bioinformatics. He has more than 100 research publications in refereed conference/workshop proceedings and journals. He was the coeditor of *Journal of Programming Languages* and a member of the editorial board of *IEEE Concurrency*. Currently, he is an associate editor of the *IEEE Transaction on Computers* and is a member of the editorial board of *Parallel Programming Letters* (PPL). He has been a program/general chair or a member of steering/program/organizing committee of many international conferences in his field. Professor Gao was a Distinguished Visitor of IEEE Computer Society from 1998-2000 and is a senior member of the IEEE, and a member of the ACM SIGARCH, SIGPLAN and IFIP WG 10.3.

► For more information on this or any computing topic, please visit our Digital Library at <http://computer.org/publications/dlib>.