

Efficient Computation of Interprocedural Control Dependence

James Ezick
Cornell University
4139 Upson Hall
Ithaca, NY 14853 USA
(607) 255-4934
ezick@cs.cornell.edu

Gianfranco Bilardi
Università di Padova
Dip. di Elettronica ed
Informatica
Via Gradenigo, 6A
35131 Padova, Italy
+390498277819
bilardi@dei.unipd.it

Keshav Pingali
Cornell University
457A Rhodes Hall
Ithaca, NY 14853 USA
(607) 255-7203
pingali@cs.cornell.edu

15 August 2001

Abstract

Control dependence information is useful for a wide range of software maintenance and testing tasks. For example, program slicers use it to determine statements and predicates that might affect the value of a particular variable at a particular program location. In the intraprocedural context an optimal algorithm is known for computing control dependence that unfortunately relies critically on the underlying intraprocedural postdominance relation being tree-structured. Hence, this algorithm is not directly applicable to the interprocedural case where the transitive reduction of the postdominance relation can be a directed acyclic graph (DAG), with nodes having multiple immediate dominators.

In this paper we present two efficient, conceptually simple algorithms for computing the interprocedural postdominance relation that can be used to compute interprocedural control dependence. For an interprocedural control flow graph $G = (V, E)$, our reachability based algorithm takes time and space $O(|V|^2 + |V||E|)$. Unlike other algorithms, it does not perform confluence operations on whole bit-vectors and can be tuned to concentrate on the interprocedural rather than intraprocedural relations in a program thus allowing it to scale better to larger programs.

Keywords

Interprocedural Analysis, Postdominance, Control Dependence, Reachability

1 Introduction

Many problems in software engineering such as testing and maintenance of programs require the computation of the *control dependence relation* of the program. Intuitively, a statement w is control dependent on a statement u in a program if there are multiple exits out of u and the choice of exit determines whether w is executed. For example, in an *if-then-else* construct statements on the two sides of the conditional statement are control dependent on the predicate.

The most common use of control dependence in software engineering is in determining whether a change to the semantics of a program statement affects the execution of another program statement. Like most program analysis problems this one is undecidable, but computing program dependences provides a reasonable approximate approach to answering such a question. The research community has focused much attention on such *slicing tools* [16, 20, 31]. For example, the *system dependence graph (SDG)* [16], an interprocedural extension of the program dependence graph [9, 23], incorporates edges for both control and data dependence to allow programs to be sliced at a point p with respect to a variable x defined or used at p .

In restructuring and optimizing compilers control dependence is used in scheduling instructions across basic-block boundaries for speculative or predicated execution [3, 10, 22], in merging program versions [15], and in automatic parallelization [2, 9, 30]. In some applications, such as code scheduling, it is necessary to know which nodes have the same control dependences as a given node. This information is useful in code scheduling because basic blocks with the same control dependences can be treated as one large basic block, as is done in region scheduling [12]. This information can also be used to decompose the control flow graph of a program into single-entry single-exit regions, and this decomposition can be exploited to speed up dataflow analysis by combining structural and fixpoint induction [17, 18] and to perform dataflow analysis in parallel [11, 18].

Formulating a precise definition of control dependence in programs with nested control structures, multi-way branches, unstructured flow of control, and procedure calls can be quite subtle. The most commonly used definition, due to Ferrante *et al.*, is based on the graph-theoretic concept of *postdominance* [9]. This work was later extended by Podgurski and Clarke who distinguished between several notions of control dependence [25]. Bilardi and Pingali [4] proposed a generalized framework to unify many different such notions. Most of the research in this area has focused on computing *intraprocedural* control dependence [9, 24, 7, 26, 27]. In this approach each procedure is treated in isolation ignoring the transfer of control due to procedure calls and returns. While adequate for some applications such as instruction scheduling, ignoring calls and returns is not an option for other applications such as interprocedural dataflow analysis. For example, a popular method to perform intraprocedural dataflow analysis consists of building *sparse dataflow representations* (such as the Static Single Assignment (SSA) form [7], sparse dataflow evaluator graphs [6], and Quick Propagation Graphs [18]), solving dataflow problems in these representations, and then projecting the solution onto the original program. Algorithms for constructing sparse representations are related to algorithms for computing control dependence in the reverse control flow graph of the program. Extending this method to the interprocedural context obviously requires the computation of interprocedural control dependence.

In the intraprocedural case, a key step toward control dependence computation is the construction of the transitive reduction of the postdominance relation. This reduced relation is tree-structured, a fact that is crucially exploited both in computing it [19, 5] and in using it for control dependence computations [24]. In the interprocedural case the transitive reduction of postdominance is not necessarily a tree, so these intraprocedural algorithms can no longer be used. The key point, to be discussed more closely in Section 2, is that while all graph-theoretic paths in the intraprocedural control flow graphs correspond to valid program executions, this is no longer true in the interprocedural case. Indeed, a procedure can be called from any number of different program points, but at the end of the call control can return only to the program point just after the point of invocation.

A search of the literature revealed only a handful of algorithms for computing interprocedural control dependence. Loyall and Mathisen [21] gave an algorithm that was improved by Harrold, Rothermel and Sinha for computing interprocedural control dependence with $O(N^5)$ complexity [13, 14]. In later work [31], they present a method for augmenting control flow graphs to handle programs with arbitrary control flow. It should be noted that their notion of control dependence differs slightly from the relation we compute in this paper in that, under their definition, the start of a procedure is always control dependent on each call site to the procedure. Their justification for this is that the call site dictates the parameters passed to the procedure and thus a natural dependence exists. It is our view that that constitutes a data dependence and thus we do not include it in our relation. This distinction does not change the complexity of the overall problem.

In this paper we present two practical algorithms for computing the interprocedural postdominance relation of a program that can be used to compute the interprocedural control dependence relation. Introduced in Section 3, one is an iterative approach and the other is based on determining reachability of nodes in the interprocedural control flow graph along so-called *valid path suffixes*. Defined in Section 2, they reflect the constraints of valid program executions. The running time of our reachability based algorithm is $O(|V|(|E| + |V|))$ where $|V|$ and $|E|$ are, respectively, the number of nodes and of edges in the interprocedural control flow graph. Precomputing and caching certain sets of reachable nodes improves the efficiency of that algorithm in practice, as we show in Section 4. Both algorithms produce the full postdominance relation, which is transitive. In Section 5 we show how, by means of a single boolean matrix multiplication, one can obtain the transitive reduction of the postdominance relation that is preferable to work with in some applications. In Section 6, we outline how to compute the interprocedural control dependence relation

of the program. We give experimental results in Section 7, comparing the performance of the reachability algorithm with that of the iterative dataflow algorithm. Finally, we discuss ongoing work in Section 8.

2 Concepts and Definitions

A program is formed by a family \mathcal{P} of procedures that can call each other, including a distinguished procedure `MAIN` that no other procedure may call. The *interprocedural* control-flow graph (ICFG) of a program models the possible transfer of control both between statements of the same procedure and between different procedures by means of the call-return mechanism [29]. As a preliminary step, it is convenient to introduce the *intraprocedural* control-flow graph (iCFG) of an individual procedure.

Definition 1 An **intraprocedural control-flow graph (iCFG)** of a procedure P is a directed graph $G_P = (V_P, E_P)$ in which nodes represent statements and an edge $u \rightarrow v$ represents possible flow of control from statement u to statement v . The node set V_P can be partitioned as

$$V_P = \{\text{START}_P\} + \{\text{END}_P\} + V_P^c + V_P^r + V_P'$$

where: START_P is a node with no predecessors from which every node is reachable; END_P is a node with no successors and reachable from every node; V_P^c is the set of call nodes, that have exactly one outgoing edge; V_P^r is the set of return nodes, that have exactly one incoming edge; $|V_P^c| = |V_P^r|$.

The edge set can be partitioned as $E_P = E_P^i \cup E_P^s$. An internal edge $(u, v) \in E_P^i$ corresponds to direct transfer of control from u to v , internally to procedure P .

The set E_P^s of short-cut edges is a subset of $V_P^c \times V_P^r$ and induces a one-to-one correspondence between call nodes and return nodes. For each $(u, v) \in E_P^s$, a $\text{label}(u, v) \in \mathcal{P}$ identifies the procedure being called at u and returning control at v . We shall use the convenient notations $v = r(u)$ and $u = c(v)$. If $\text{label}(u, v) = F$, we shall say that u is a call node for F and v is a return node from F .

The ICFG is obtained by assembling the iCFGs of all procedures and replacing each short-cut edge (u, v) having $\text{label}(u, v) = F$ with two interprocedural edges (u, START_F) and (END_F, v) , that are said to correspond to each other.

Definition 2 The **interprocedural CFG (ICFG)** of a program is a graph $G = (V, E)$ with $V = \cup_{P \in \mathcal{P}} V_P$ and $E = E^i + E^c + E^r$, where $E^i = \cup_{P \in \mathcal{P}} E_P^i$ is called the set of internal edges, and E^c and E^r , defined as

$$E^c = \cup_{P \in \mathcal{P}} \{(u, \text{START}_F) : (u, v) \in E_P^s, \text{label}(u, v) = F\},$$

$$E^r = \cup_{P \in \mathcal{P}} \{(\text{END}_F, v) : (u, v) \in E_P^s, \text{label}(u, v) = F\}$$

and called the sets of the call edges and of the return edges, respectively.

Figure 1(a) shows an interprocedural CFG for a simple program with two procedures `MAIN` and `F` in which `F` is called from two places in `MAIN`.

Not every ICFG path corresponds to a possible path of execution since a procedure can be invoked from many call sites but each invocation can only return to its corresponding call site. For example, in Figure 1(a), the path $(a, b, c, j, k, l, f, h, i)$ does not correspond to a valid path of execution since the call to `F` at node c does not return to node e . This intuition is captured by defining a *valid path* (this is called a *complete* valid path by Sharir and Pnueli [29]) as a path where the subsequence of call and return edges is *proper*, in the following sense.

Definition 3 The set of **proper sequences** of call and return edges is defined as the context-free language on the alphabet $(E^c + E^r)$ that (i) contains the empty sequence, (ii) is closed under concatenation, and (iii) if it contains sequence σ then it also contains those sequences of the form $(u, \text{START}_F)\sigma(\text{END}_F, v)$, for each short-cut edge (u, v) with $\text{label}(u, v) = F$.

Definition 4 A path from $\text{START}_{\text{MAIN}}$ to END_{MAIN} in the ICFG G of a program, viewed as a sequence of edges, is a **valid path** if its subsequence of call and return edges is proper.

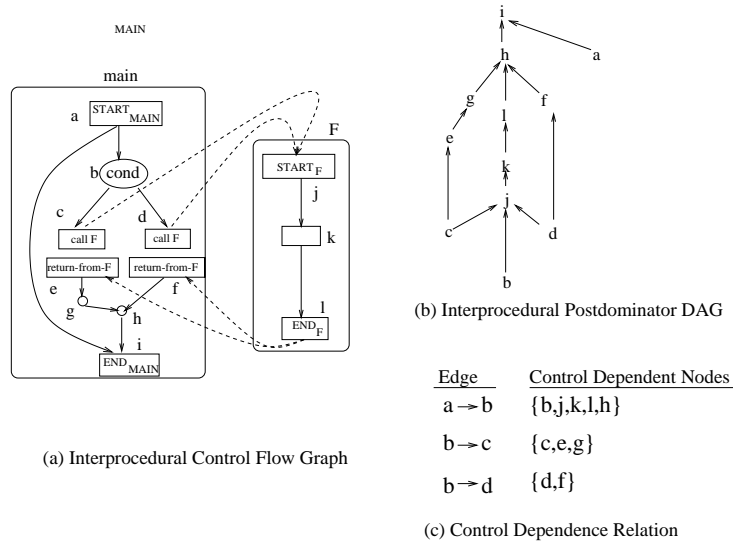


Figure 1: A Small Program, and its Postdominator DAG and Control Dependence Relation

In the remainder of this paper, we shall restrict our attention to ICFGs where each node v occurs on some valid path. We will have occasion to deal with the following kind of valid path segment, referred to as *same-level* valid path segments in the literature [28].

Definition 5 A *same-level valid path segment* is a sequence of edges that appear consecutively in some valid path, whose first and last nodes belong to the same procedure, and whose subsequence of call and return edges is proper.

In our running example, path (b, c, j, k, l, e) is a same-level valid path segment.

The following lemma is needed in the proof of correctness of the algorithm for computing interprocedural postdominance.

Lemma 1 Let σ_1 and σ_2 be two same-level valid path segments from node u to node v . If $\pi = \pi_1 \sigma_1 \pi_2$ is a valid path, then so is $\tau = \pi_1 \sigma_2 \pi_2$.

Proof: It is easy to see that τ is a path from the start of the program to its end. By Definition 5, the subsequences of call and return edges in both σ_1 and σ_2 are proper sequences. Furthermore, given a proper sequence of call and return edges, we can always replace a proper subsequence with another proper subsequence to obtain a proper sequence. Therefore, τ is a valid path. \square

We can now define interprocedural dominance and postdominance between nodes of the ICFG.

Definition 6 A node v is said to **interprocedurally dominate** a node w if v occurs before w in every valid path that contains w .

A node v is said to **interprocedurally postdominate** a node w if v occurs after w in every valid path that contains w .

It is easy to show that dominance and postdominance are transitive relations. Figure 1(b) shows the transitive reduction of the postdominance relation of the program of Figure 1(a).

Control dependence can be introduced formally as follows [9]:

Definition 7 A node $w \in V$ is said to be **control dependent** on edge $(u \rightarrow v) \in E$ if

1. w postdominates v , and
2. if $w \neq u$, then w does not postdominate u .

Call Nodes: $PDOM(\text{Call}_F) = PDOM(\text{START}_F) \cup PDOM(\text{Return}_F) \cup \{\text{Call}_F\}$

END_{MAIN}: $PDOM(\text{END}_{\text{MAIN}}) = \{\text{END}_{\text{MAIN}}\}$

Other Nodes: $PDOM(n) = \bigcap_{s \in \text{succ}(n)} PDOM(s) \cup \{n\}$

Figure 2: Rules for postdominance dataflow equations

By convention, nothing is control dependent on a return edge since the END of a procedure is not a decision point although there may be multiple edges emanating from it.

Intuitively, if control flows from node u to node v along edge $u \rightarrow v$, it will eventually reach node w ; however, control may reach END from u without passing through w . Thus, u is a ‘decision-point’ that influences the execution of w .

Definition 8 *Given an ICFG $G = (V, E)$, its **control dependence relation** is the set $C \subseteq E \times V$ of all pairs (e, w) such that node w is control dependent on edge e .*

The control dependence relation for the program of Figure 1(a) is shown in Figure 1(c).

Typically, both software engineering and compiler applications of control dependence require the computation of the following sets derived from C [8]:

Definition 9 *Given a node w and an edge e in an ICFG with control dependence relation C , we define the following **control dependence sets**:*

- $\text{cd}(e) = \{w \in V : (e, w) \in C\}$,
- $\text{conds}(w) = \{e \in E : (e, w) \in C\}$, and
- $\text{cdequiv}(w) = \{v \in V : \text{conds}(v) = \text{conds}(w)\}$.

3 Two Algorithms for Computing Interprocedural Postdominance

In this section, we present two algorithms for computing the interprocedural postdominance relation. The first algorithm is a relatively straight-forward dataflow algorithm that solves a set of monotone dataflow equations iteratively. The second algorithm is based on computing graph-reachability in subgraphs of the interprocedural CFG.

3.1 Iterative Algorithm

It is well-known that the intraprocedural postdominance relation can be computed by solving a set of monotone dataflow equations [19]. This is true for the interprocedural postdominance relation as well. The lattice underlying the equations is the powerset of nodes in the program, ordered by containment, in which the least element is the empty set. The postdominance relation can be described implicitly by writing down a set of equations, one per node, in which the postdominators of a node are expressed as a function of the postdominators of its successors in the interprocedural CFG and some other nodes. Figure 2 shows the rules for generating these equations. In general, the postdominators of a node n are n itself and any other node that postdominates all the control flow successors of n . A Call node is also postdominated by any node that postdominates its corresponding Return node.

This is a monotone dataflow problem. Its solution can be found in the usual way by iterating downward from the initial approximation that assumes the postdominator set of each node is the set of all program nodes.

As is usual, our implementation maintains a work-list of nodes whose postdominator set must be recomputed because one or more of the sets on the right hand side of its equation has changed. This work-list is initialized to END_{MAIN}, and nodes are enqueued and dequeued till convergence occurs. Postdominator sets are represented as bit-vectors, and the necessary union and intersection operations are performed using bit-vector operations.

3.2 Reachability-based Algorithm

Let us observe that, by a straightforward reformulation of Definition 6, a node v fails to postdominate precisely those nodes w for which there is a valid path suffix $w \rightarrow \text{END}_{\text{MAIN}}$ that does not contain v . This observation is the basis for our algorithm, the pseudocode for which is shown in Figure 3. This code assumes that each node is given a unique number between 1 and the total number N of nodes in the ICFG, and that corresponding call and return nodes can be determined from each other in constant time.

The algorithm determines the interprocedural postdominance relation PDOM by first initializing it to contain all pairs (v, w) (line 3) and then pruning from it, for each node v (processed by one iteration of the loop in line 7), all those nodes w that are reachable from END_{MAIN} in the reverse ICFG with v deleted, via the reverse of a valid path suffix. The actual pruning is accomplished by the call `SEARCH-WITHOUT(v)`. This procedure (line 10) determines reachability by performing a search. A work-list of reachable nodes is maintained, initialized to contain END_{MAIN} (line 13), and from which nodes are extracted one at the time (line 15). A node w is marked when it is first encountered (line 17), to avoid processing it multiple times, and the entry (v, w) is deleted from PDOM (line 19). Then, procedure `VISIT` is invoked (line 20) to add the predecessors of w to the work-list, where appropriate.

Reachability along valid path segments introduces some subtleties, specifically when $w = \text{START}_F$, for some procedure F . In this case, a generic predecessor c of w is a call node; hence a valid path suffix from w to END_{MAIN} can be extended by prepending edge (c, w) if and only if such a path includes the corresponding return node $r(c)$. Care must be taken to handle correctly and efficiently the situation when START_F is visited before $r(c)$ as well as the situation when the visits happen in the opposite order. To this end, when a return node $r(c)$ corresponding to a call node c for F is visited, in addition to adding END_F to the work-list (line 26) node c is considered as well: if START_F has already been visited, then c is added to the work-list (line 28), else c is inserted (line 29) into an initially empty *bucket* associated with procedure F . When START_F is visited, the bucket - which contains call nodes whose corresponding return has already been visited - is simply emptied into the work-list (lines 30-35). For a node w that is neither a return nor a start node, the visit simply adds all predecessors to the work-list (lines 36-40).

Deletion of node v from the ICFG is not actually performed, but simulated simply by skipping the visit of v (line 18) and hence preventing the search from propagating beyond v .

Next, we establish the correctness and analyze the performance of our algorithm.

Theorem 1 *Given as input an ICFG $G = (V, E)$, procedure `COMPUTE-PDOM` runs in time $O(|V|(|V| + |E|))$ and sets $\text{PDOM}[v][w]$ to true if and only if v interprocedurally postdominates w .*

Proof: The performance bound simply follows from the fact that $|V|$ searches are executed, each taking time proportional to the number $|V|$ of nodes and $|E|$ of edges.

Correctness requires a more detailed argument. We first consider separately the case $v = \text{END}_{\text{MAIN}}$, where v is added to the work-list in line 13 and is removed from the work-list in line 15. The while loop in line 14 executes just once and, for each w , the entry $\text{PDOM}[\text{END}_{\text{MAIN}}][w]$ will remain true, which is correct since END_{MAIN} postdominates all nodes.

Now, let $v \neq \text{END}_{\text{MAIN}}$. We will show that $\text{PDOM}[v][w]$ will be set to false iff there is a non-empty valid path suffix from w to END_{MAIN} that does not contain v .

Notationally, $y \rightarrow z$, $y \xrightarrow{*} z$, and $y \xrightarrow{\pm} z$ respectively denote an edge, a path, and a non empty path between nodes y and z .

- \Leftarrow : We show inductively that, for every n , if $\pi = w \xrightarrow{\pm} \text{END}_{\text{MAIN}}$ is a non-empty valid path suffix of length n that does not contain v , then $\text{PDOM}[v][w]$ is set to false.

If $n = 0$, then $w = \text{END}_{\text{MAIN}}$. This node is added to the work-list in line 13, removed from it in line 15, and $\text{PDOM}[v][w]$ is set to false in line 19.

Assume now that the inductive hypothesis holds for lengths no larger than n . Now consider a node w for which there is a non-empty valid path suffix $w \rightarrow x \xrightarrow{*} \text{END}_{\text{MAIN}}$ of length $n+1$. By the inductive assumption, $\text{PDOM}[v][x]$ will be set to false at some point. This must happen at line 19, which is immediately followed by line 20 which invokes procedure `VISIT`. Consider the three possible cases for node x .

1. If node x is a return node and w is not marked, then w is added to the work-list in line 26. At some later point, it is removed from the work-list and $\text{PDOM}[v][w]$ is set to false. On the other hand, if w is marked, it must have been marked in line 17 which is followed by line 19 in which $\text{PDOM}[v][w]$ is set to false. In either case, the required result follows.

```

1. ICFG G;
2. int N = number of nodes in G;
3. boolean[1..N][1..N] PDOM = true;
   // PDOM[v][w] = false if v does not postdominate w
4. int P = number of procedures in G;
5. node_set work-list = {}; //set of nodes to be visited by graph search

6. procedure COMPUTE-PDOM () {
7.   FOR v = 1..N DO
8.     // find nodes reachable from END-MAIN in reverse ICFG w/o v
9.     SEARCH-WITHOUT(v);
10.  END
11. }

12. procedure SEARCH-WITHOUT(node v) {
13.  node_set[P] buckets = {};
14.  boolean[N] marked = false; // no node initially visited

15.  work-list.add(END-MAIN);
16.  // begin reverse reachability at END-MAIN
17.  WHILE(work-list <> empty) DO
18.    node w = work-list.remove();
19.    IF (marked[w]) THEN continue; // w has already been visited
20.    marked[w] = true;
21.    IF (w == v) continue; // skip v since it is conceptually removed from G
22.    PDOM[v][w] = false; // w is reachable from END-MAIN w/o v
23.    VISIT(w, v, buckets, marked); // process predecessors of w
24.  END
25. }

26. procedure VISIT(node w, node v, node_set[] buckets, boolean[] marked) {
27.  // process predecessors of w
28.  IF (w is a return node from procedure F) // w target of return edge
29.  THEN
30.    IF (not marked(END-F)) THEN
31.      work-list.add(END-F);
32.      let node c = c(w) // call node corresponding to w;
33.      IF (marked[START-F]) THEN work-list.add(c);
34.      ELSE buckets[F].add(c);
35.    ELSEIF (w is a START node for procedure F) // w target of call edge
36.    THEN
37.      FOR each node c in buckets[F] DO
38.        remove c from buckets[F];
39.        work-list.add(c);
40.      END
41.    ELSE
42.      FOR each predecessor z of w DO
43.        IF (not marked[z]) THEN
44.          work-list.add(z);
45.        END
46.      END
47.  }

```

Figure 3: Algorithm for computing the interprocedural postdominance relation

2. If node x is the **START** node for some procedure F , w must be a call node for F and the corresponding return node $r(w)$ must occur on the valid path suffix $x \xrightarrow{*} \text{END}_{\text{MAIN}}$. By the inductive assumption, both x and $r(w)$ must be put on the work-list, and procedure **VISIT** must be called with both these nodes. If x is processed first, then line 28 is executed when $r(w)$ is processed at a later time, and w is put on the work-list. On the other hand, if $r(w)$ is processed first, then w is added to $\text{buckets}[F]$ in line 29, and w is added to the work-list in line 34 when x is processed. In either case, at some point, w is extracted from the work list and $\text{PDOM}[v][w]$ is set to false.
 3. The remaining case is when $w \rightarrow x$ is an internal edge. When x is processed, w is either already marked, in which case $\text{PDOM}[v][w]$ is set to false, or w is not marked, in which case it is added to the work-list, and $\text{PDOM}[v][w]$ is set to false when w is processed.
- \Rightarrow : We note that elements of the PDOM array are set to false in the while loop of line 14 and we show inductively that, for each k , if $\text{PDOM}[v][w]$ is set to false at the k -th iteration, then there is a valid path suffix π from w to END_{MAIN} that does not contain v .

The while loop executes at least once since node END_{MAIN} is added to the work-list in line 13. $\text{PDOM}[v][\text{END}_{\text{MAIN}}]$ is set to false in the first iteration, and the required result is trivially obtained by letting π be the empty path, obviously a valid path suffix.

Assume now that the inductive assumption holds for the first k iterations of the while loop. Let w be the node removed from the work-list in iteration $k + 1$ and assume $w \neq v$. Therefore, node w is distinct from END_{MAIN} and it must have been added to the work-list when some node r was processed in procedure **VISIT**. By the inductive assumption, there is a valid path

$$\pi = \text{START}_{\text{MAIN}} \xrightarrow{*} r \xrightarrow{+} \text{END}_{\text{MAIN}}$$

and the suffix $r \xrightarrow{+} \text{END}_{\text{MAIN}}$ does not contain v . Consider the three possible cases for node r .

1. Suppose r is a return node for procedure F .
Suppose $w = \text{END}_F$. Since the only ICFG predecessor of node r is w , the valid path π must contain w , so there is a valid path suffix $w \xrightarrow{+} \text{END}_{\text{MAIN}}$ that does not contain v .
Otherwise, $w = c(r)$ is the call node corresponding to r . Then, any valid path that contains r must contain w since r so π contains w , and there is a valid path suffix

$$\rho = w \rightarrow \text{START}_F \xrightarrow{+} \text{END}_F \xrightarrow{+} \text{END}_{\text{MAIN}}.$$

By the inductive assumption, there is a valid path suffix $\sigma = \text{START}_F \xrightarrow{+} \text{END}_F \xrightarrow{+} \text{END}_{\text{MAIN}}$ that does not contain v . Path σ has a same-level valid path segment $\gamma = \text{START}_F \xrightarrow{+} \text{END}_F$. Replacing the same-level valid path segment from START_F to END_F in ρ with γ , we get a valid path suffix $w \xrightarrow{+} \text{END}_{\text{MAIN}}$ that does not contain v .

2. Suppose r is a **START** node for procedure F .
Then w must be a call node for F which was in $\text{buckets}[F]$. This node must have been added to $\text{buckets}[F]$ in line 29 when the return node r' corresponding to w was processed by procedure **VISIT**. Node r' must have been processed before node r , so by inductive assumption, there is a valid path $\text{START}_{\text{MAIN}} \xrightarrow{*} r' \xrightarrow{+} \text{END}_{\text{MAIN}}$ which must contain w , so the required result holds.
3. The final case is that $w \rightarrow r$ is an internal edge.
By assumption, there is a valid path
 $\pi = \text{START}_{\text{MAIN}} \xrightarrow{*} r \xrightarrow{+} \text{END}_{\text{MAIN}}$. Let F be the procedure in which r occurs. By assumption, START_F must occur on π , which can then be written as

$$\pi = \text{START}_{\text{MAIN}} \xrightarrow{*} \text{START}_F \xrightarrow{*} r \xrightarrow{+} \text{END}_{\text{MAIN}}$$

By assumption about programs, there is a same-level valid path segment $\beta = \text{START}_F \xrightarrow{*} w \rightarrow r$. Replacing the segment $\text{START}_F \xrightarrow{*} r$ with β , we get a valid path that contains w as required; furthermore, the suffix $w \xrightarrow{+} \text{END}_{\text{MAIN}}$ does not contain v .

□

4 Precomputing Reachability

We now show that the efficiency of the algorithm for computing the interprocedural postdominance relation can be improved by precomputing and caching *intraprocedural* reachability information and using this

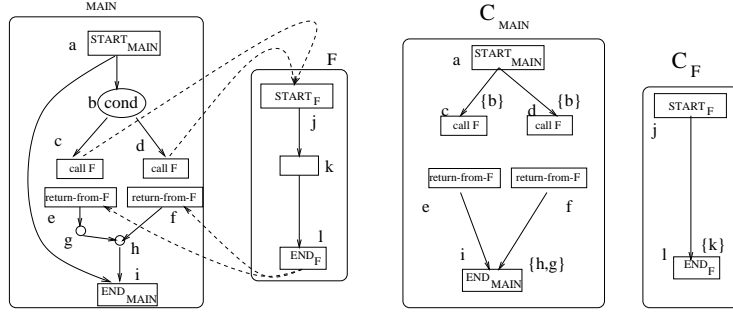


Figure 4: Collapsed Graphs for Running Example

information selectively to ensure that most of the graph traversals performed during the postdominance computation are along interprocedural edges.

In the preprocessing step, for each procedure $P \in \mathcal{P}$, we perform reachability computations in the intraprocedural control flow subgraph $S_P = (V_P, E_P^0)$, where the short-cut edges have been removed. The nodes in $V_P^c \cup \{\text{END}_P\}$ (namely, the call nodes and END_P) are called *root nodes* of this graph. We introduce the *locally reachable node set* $L(r)$ that is the set of nodes in V_P^l (nodes that are not START_P , END_P or call or return nodes, as in Definition 1) reachable from root node r in the reverse graph of S_P . Intuitively, in the reverse ICFG, nodes in $L(r)$ can be reached from r without traversing interprocedural edges. We then build a collapsed graph $C_P = (I_P, D_P)$ in which the nodes are START_P , END_P , and all the call and return nodes of P , and in which there is an edge (m, n) if $m, n \in I_P$ and $m \in L(n)$. At the END_P node and each call node in this collapsed graph we store the locally reachable set of nodes computed for the corresponding node in S_P (the algorithm below assumes that this set is stored in an array named L , indexed by the node). During the postdominator computation, we use S_P for reachability computations if the deleted node is not in P ; otherwise, we use the original ICFG for P .

Figure 4 shows the collapsed graphs for the procedures in the running example. Figure 5 shows the modifications that are required to the algorithm of Figure 3 when preprocessing is used.

The proof of correctness of the algorithm with preprocessing hinges on the fact that the intraprocedural reachability computation that determines locally reachable sets finds same-level paths within procedures that, prepended to paths from the END nodes of the procedure to END_{MAIN} , yield valid path suffixes.

5 Computation of Transitive Reduction

PDOM is a reflexive, anti-symmetric, transitively closed relation. Boolean array $\text{PDOM}[v][w]$ can be viewed as the adjacency matrix of a graph where there is an edge $(u \rightarrow v)$ in the graph if u postdominates v . This graph has self-loops at every node (reflexivity); except for self-loops, it is acyclic (anti-symmetry) and, if $(u \rightarrow v)$ and $(v \rightarrow w)$ are edges in the graph, so is edge $(u \rightarrow w)$ (transitive closure).

For some applications, it is useful to work the transitive reduction IPDOM , which can be computed by a single boolean matrix multiplication, as we show in this section.

The first step is to define the irreflexive version P_I of PDOM , by the equation

$$P_I = \text{PDOM} \wedge \neg I_n, \quad (1)$$

where I_n is the $n \times n$ identity matrix, whose diagonal entries are true and the other are false. We show next that the transitive reduction IPDOM of PDOM can be computed by the following expression:

$$\text{IPDOM} = P_I \wedge \neg P_I^2. \quad (2)$$

Lemma 2 *The relation IPDOM is transitively reduced.*

Proof: Note first that if $(i, k) \in \text{IPDOM}$, then $(i, k) \in P_I$ (from Equation 2), so $(i, k) \in \text{PDOM}$ and i is distinct from k (from Equation 1).

```

// Assume that CF is the collapsed graph for procedure F
.....
//reverse reachability from END-SMAIN if v belongs to MAIN, END-CMAIN otherwise
19a. IF (v and w do not belong to the same procedure) && (w is a call or end node)
    b. THEN
    c.   FOR each node n in L[w] DO
    d.     PDOM[v][n] = false;
    e.   ELSE PDOM[v][w] = false;
.....
22. procedure VISIT(node w, node v, node_set[] buckets, boolean[] marked) {
    // process predecessors of w
23.   IF (w is a return node in G from procedure F) // w target of return edge
24a.  THEN
    b.    IF (v belongs to F) THEN let GRAPH-F = SF // use full graph of F
    c.    ELSE let GRAPH-F = CF; // use collapsed graph of F
25.    IF (not marked(END-GRAPH-F)) THEN
26.      work-list.add(END-GRAPH-F);
27.    let node c = c(w) // call node paired with w in graph (SG or CG) containing w
28.    IF (marked[START-GRAPH-F]) THEN work-list.add(c);
29.    ELSE buckets[F].add(c);
.....
    }

```

Figure 5: Postdominator Computation with Preprocessing

Suppose $(i, j), (j, k) \in IPDOM$. We have just shown that $(i, j), (j, k) \in P_I$, so $(i, k) \in P_I^2$. Therefore, $(i, k) \notin IPDOM$. \square

Theorem 2 $IPDOM^* = PDOM$.

Proof:

- $IPDOM^* \subseteq PDOM$: We show that $IPDOM \subseteq PDOM$; since $PDOM$ is transitively closed this proves the required result. From the definitions of $IPDOM$ and P_I , we obtain the following equation:

$$IPDOM = P_I \wedge \neg P_I^2 = PDOM \wedge \neg I_n \wedge P_I^2$$

The conjunction of $PDOM$ with other matrices never has more *true* entries than $PDOM$ itself, so $IPDOM \subseteq PDOM$.

- $PDOM \subseteq IPDOM^*$: If $(x, y) \in PDOM$, let

$$x = z_0, z_1, z_2, \dots, z_n = y$$

be the longest path from x to y in the graph of $PDOM$ where z_0, z_1, \dots, z_n are all distinct nodes (such a path exists because the graph, without self-loops, is acyclic). Let us call this path R .

For all k , z_k and z_{k+1} are distinct nodes, so $(z_k, z_{k+1}) \in P_I$; furthermore, $(z_k, z_{k+1}) \notin P_I^2$ since otherwise, there is a path of length 2 from z_k to z_{k+1} in the graph of P_I , and hence in the graph of $PDOM$, contradicting the assumption that R is the longest path from x to y in the graph of $PDOM$. \square

These facts lead to the following result.

Lemma 3 *The transitive reduction of the interprocedural postdominance relation can be computed in time $O(|V| * (|E| + |V|^2))$.*

6 Computation of Interprocedural Control Dependence

In this section, we briefly outline how to compute the control dependence relation C and its related sets (Definitions 8 and 9). One can represent C as an $|E| \times |V|$ Boolean array whose entry $C[e, w]$ is set to true

| Code | V | E | P | V^c | PDOM | I. Set | I. List | Time (ms) | R. Set | R. List | Time |
|-------------|-----|-----|-----|-------|-------|---------|---------|-----------|--------|---------|-------|
| mutual | 35 | 40 | 5 | 7 | 341 | 6098 | 155 | 470 | 1297 | 1454 | 440 |
| postfix | 164 | 209 | 8 | 26 | 3736 | 153552 | 803 | 2020 | 28630 | 34463 | 1960 |
| tic-tac-toe | 229 | 276 | 13 | 24 | 9842 | 326713 | 1263 | 2470 | 37017 | 40710 | 2910 |
| compress | 585 | 684 | 21 | 37 | 21489 | 2399114 | 3608 | 10570 | 189304 | 210340 | 14800 |

Table 1: Algorithm Performance Results

iff node w is control dependent on edge e . Alternatively, C can be considered as the edge set of a bipartite graph $B = (E, V, C)$, with vertex sets E and V . Both representations have size $O(|V||E|)$, in the worst case; if the relation is sparse, however, the size of B could be much smaller.

Relation C can be easily constructed from the interprocedural postdominance relation by simply scanning all pairs $(e = (u, v), w)$ and checking the entries $\text{PDOM}[w][u]$ and $\text{PDOM}[w][v]$. This takes time $O(|V||E|)$, for both the array and the graph representations.

In the graph representation, set $\text{cd}(e)$ is easily obtained, in time proportional to its size, by collecting the neighbors of e in B . A similar approach can be used to obtain $\text{conds}(w)$.

Finally, the cdequiv sets, that are the equivalence classes of nodes that have the same dependences, can be obtained by the following approach. A partition of V , initially consisting of just one set, is progressively refined until, at the end, the blocks of the partition coincide with the cdequiv sets. A refinement phase is performed for each edge e by splitting each block of the partition into two subblocks, respectively containing the nodes that are and are not control dependent on e . It is easy to see that the entire process can be completed in time $O(|V||E|)$.

7 Experimental Results

We have implemented these algorithms as a plug-in to the GrammaTech CodeSurfer [1] program analysis tool. We used only the control flow graphs generated by CodeSurfer for C programs and none of that tool's other analysis capabilities. The implementation of each algorithm (iterative and reachability with preprocessing) required less than 100 lines of scheme code.

Table 1 compares the performance of the two algorithms for computing interprocedural postdominance that we have discussed in this paper. We used the following test programs:

- mutual: A test program involving mutual recursion
- postfix: A postfix calculator
- tic-tac-toe: A text based tic-tac-toe game using alpha-beta search
- compress: SPECint95 file compression utility

The first four data columns in Table 1 give the number of nodes, edges, procedures, and call nodes in the program respectively. The column $PDOM$ shows the number of non-zero entries in the full interprocedural postdominance relation. The columns $I. Set$ and $I. List$ indicate the number of set and work-list operations performed by the iterative algorithm respectively. Likewise, the columns $R. Set$ and $R. List$ indicate the number of set and work-list operations performed by the reachability algorithm. The time columns are in milliseconds as recorded on our 650Mhz Pentium III (256 MB RAM) system.

Since the reachability algorithm makes only local updates to the $PDOM$ relation when a node is processed it performs far fewer set update operations than the iterative algorithm. The iterative algorithm must perform union and intersection operations over bit vectors of length equal to the number of nodes in the program at each update. By storing the results of our preprocessing searches in small, single-procedure-sized bit-vectors, we are able to get most of the advantage of updating a word length bit vector segment in a single operation without the added expense of updating an entire $|V|$ -sized bit vector. Furthermore, since the reachability based algorithm relies upon repeated searches, it generally requires more constant-time list operations than the iterative algorithm. Therefore, in practice, the two algorithms offer competitive performance.

It should be noted that for very large programs, the reachability based algorithm, which computes a single column of the postdominance relation in each iteration, should demonstrate significantly better data

locality than the iterative algorithm. The reachability based algorithm also offers the ability to compute only the nodes a particular node postdominates; a feature the iterative algorithm does not offer.

8 Conclusions and Future Work

In this paper we have presented two efficient algorithms for computing interprocedural control dependence by first computing interprocedural postdominance. While the reachability based algorithm we present runs in time quadratic in the size of the CFG, it requires that the entire postdominance relation be computed. Our future goal is to develop an algorithm to compute the postdominance DAG directly in a manner similar to the one by which Lengauer and Tarjan [19] directly compute the postdominance tree in the intraprocedural case. Such an algorithm would eliminate the need to compute a transitive reduction as well as ameliorate the space demands of present algorithms that often make them infeasible for very large programs.

From the DAG representation of the postdominance relation our hope is to extend the Roman chariots [24] formulation of control dependence queries to this case. Such an extension would allow us to answer control dependence queries more quickly and without the need to either store or explicitly compute the entire control dependence relation.

It is also our goal to incorporate into our program model atypical control flow effects such as embedded halts and exception handling. Recent work [31] suggests that such effects can be incorporated by augmenting the existing program representation with additional nodes and edges and then acting in some appropriate way when traversing the additional edges to preserve context. It is our belief that the core algorithm presented here will extend naturally to those representations.

Finally, we hope to integrate these algorithms into a toolkit to do both interprocedural control and data flow analysis as well as program slicing on large programs. The accepted data structure for slicing is the system dependence graph, *SDG*, [16] that can be thought of as a multi-entry, multi-exit variant of the interprocedural control flow graph. In the *SDG* call sites and entry nodes are augmented with nodes representing actual and formal input data, respectively. Likewise the return and exit nodes are augmented with actual and formal output nodes. Hence, the formal parameter nodes act as additional entry and exit points of a procedure. Edges in the *SDG* represent traditional types of control and data dependence. Given this representation program slicing can be reduced to a reachability problem from a designated starting point. We believe that the reachability-based algorithm described in this paper extends naturally to slicing in this representation.

9 Acknowledgments

The work of G. Bilardi was supported in part by MURST of Italy and by IBM while this author was a visitor at the T. J. Watson Research Laboratory. J. Ezick and K. Pingali were supported by NSF grants EIA-9726388, ACI-9870687, EIA-9972853, ACI-0085969, ACI-0090217, and ACI-0121401.

References

- [1] www.grammotech.com/products/codesurfer/codesurfer_index.html.
- [2] F. Allen, M. Burke, R. Cytron, J. Ferrante, W. Hsieh, and V. Sarkar. A framework for determining useful parallelism. In *Proceedings of the 1988 International Conference on Supercomputing*, pages 207–215, St. Malo, France, July 4–8, 1988.
- [3] D. Bernstein and M. Rodeh. Global instruction scheduling for superscalar machines. In *Proceedings of the SIGPLAN '91 Conference on Programming Language Design and Implementation*, pages 241–255, Toronto, Ontario, June 26–28, 1991.
- [4] G. Bilardi and K. Pingali. A framework for generalized control dependence. In *Proceedings of the ACM SIGPLAN '96 conference on Programming language design and implementation*, pages 291–300, May 1996.

- [5] A. L. Buchsbaum, H. Kaplan, A. Rogers, and J. R. Westbrook. Linear-time pointer-machine algorithms for least common ancestors, mst verification, and dominators. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*, May 1998.
- [6] J.-D. Choi, R. Cytron, and J. Ferrante. Automatic construction of sparse data flow evaluation graphs. In *Proceedings of the 18th ACM Symposium on Principles of Programming Languages*, pages 55–66, January 1991.
- [7] R. Cytron, J. Ferrante, B. K. Rosen, M. N. Wegman, and F. K. Zadeck. Efficiently computing static single assignment form and the control dependence graph. *ACM Transactions on Programming Languages and Systems*, 13(4):451–490, October 1991.
- [8] R. Cytron, J. Ferrante, and V. Sarkar. Compact representations for control dependence. In *Proceedings of the SIGPLAN '90 Conference on Programming Language Design and Implementation*, pages 337–351, White Plains, New York, June 20–22, 1990.
- [9] J. Ferrante, K. J. Ottenstein, and J. D. Warren. The program dependence graph and its use in optimization. *ACM Transactions on Programming Languages and Systems*, 9(3):319–349, July 1987.
- [10] J. Fisher. Trace scheduling: a technique for global microcode compaction. *IEEE Transactions on Computers*, 7(3):478–490, 1981.
- [11] R. Gupta, L. Pollock, and M. L. Soffa. Parallelizing data flow analysis. In *Proceedings of the Workshop on Parallel Compilation*, Kingston, Ontario, May 6–8, 1990. Queen’s University.
- [12] R. Gupta and M. L. Soffa. Region scheduling. In *2nd International Conference on Supercomputing*, pages 141–148, 1987.
- [13] M. J. Harrold, G. Rothermel, and S. Sinha. Computation of interprocedural control dependence. Technical Report OSU-CISRC-7/97-TR36, The Ohio State University, July 1997.
- [14] M. J. Harrold, G. Rothermel, and S. Sinha. Computation of interprocedural control dependence. In *Software Engineering Notes, Proceedings of the International Symposium on Software Testing and Analysis*, volume 23, pages 11–20. ACM/SIGSOFT, March 1998.
- [15] S. Horwitz, J. Prins, and T. Reps. Integrating non-interfering versions of programs. In *Conference Record of the 14th Annual ACM Symposium on Principles of Programming Languages*, pages 133–145, Munich, West Germany, Jan. 21–23, 1987.
- [16] S. Horwitz, T. Reps, and D. Binkley. Interprocedural slicing using dependence graphs. *Transactions on Programming Languages and Systems*, 12(1):22–60, January 1990.
- [17] R. Johnson. *Efficient Program Analysis using Dependence Flow Graphs*. PhD thesis, Cornell University, August 1994.
- [18] R. Johnson, D. Pearson, and K. Pingali. The program structure tree: Computing control regions in linear time. In *Proceedings of the SIGPLAN '94 Conference on Programming Language Design and Implementation*, pages 171–185, Orlando, Florida, June 20–24, 1994.
- [19] T. Lengauer and R. E. Tarjan. A fast algorithm for finding dominators in a flowgraph. *ACM Transactions on Programming Languages and Systems*, 1(1):121–141, July 1979.
- [20] D. Liang and M. J. Harrold. Slicing objects using system dependence graph. In *International Conference on Software Maintenance*, November 1998.
- [21] J. Loyall and S. Mathisen. Using dependence analysis to support the software maintenance process. In *Proceedings of the Conference on Software Maintenance*, pages 282–291, September 1993.
- [22] C. Newburn, D. Noonburg, and J. Shen. A PDG-based tool and its use in analyzing program control dependences. In *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques*, August 1994.

- [23] K. J. Ottenstein and L. M. Ottenstein. The program dependence graph in a software development environment. In *Proceedings of the ACM SIGSOFT/SIGPLAN Software Engineering Symposium on Practical Software Development Environments*, pages 177–184, April 1984.
- [24] K. Pingali and G. Bilardi. Optimal control dependence and the roman chariots problem. *ACM Transactions on Programming Languages and Systems*, 19(3):1–30, May 1997.
- [25] A. Podgurski and L. A. Clarke. A formal model of program dependences and its implications for software testing, debugging, and maintenance. *IEEE Transactions on Software Engineering*, 16(9):965–79, September 1990.
- [26] L. Pollock and M. Soffa. An incremental version of iterative dataflow analysis. *IEEE Transactions on Software Engineering*, 15(12):1537–1549, December 1989.
- [27] B. G. Ryder and M. C. Paull. Incremental data-flow analysis algorithms. *ACM Transactions on Programming Languages and Systems*, 10(1):1–50, January 1988.
- [28] M. Sagiv, T. Reps, and S. Horwitz. Precise interprocedural data flow analysis for with applications to constant propagation. Technical Report TR-1284, Computer Science Department, University of Wisconsin, Madison, August 1995.
- [29] M. Sharir and A. Pnueli. Two approaches to interprocedural data flow analysis. In S. Muchnick and N. Jones, editors, *Program Flow Analysis: Theory and Applications*, chapter 7, pages 189–233. Prentice-Hall, 1981.
- [30] B. Simons, D. Alpern, and J. Ferrante. A foundation for sequentializing parallel code. In *Conference Record of SPAA '90: ACM Symposium on Parallel Algorithms and Architecture*, July 1990.
- [31] S. Sinha, M. J. Harrold, and G. Rothermel. System-dependence-graph-based slicing of programs with arbitrary interprocedural control flow. In *International Conference on Software Engineering*, pages 432–441, May 1999.