

Interprocedural Analysis and Optimization

KEITH D. COOPER

MARY HALL

KEN KENNEDY

LINDA TORCZON

*Rice University
Houston, Texas*

1. Introduction

As machines and languages have become more complicated, compiler technology has necessarily become more sophisticated. With the advent of vector and parallel computers, single-procedure analysis will no longer be enough to produce high-quality parallelized code. In this paper, we introduce some of the problems that can be solved with the help of an interprocedural analysis and optimization system. We also present methods for the solution of these problems, summarizing research on interprocedural compilation over the past 15 years.

We begin with two definitions. *Interprocedural analysis* refers to gathering information about the entire program instead of a single procedure. This is analogous to the information collected by a single procedure data-flow analyzer. Examples of interprocedural analysis problems include determining the variables that might be modified as side effect of a procedure call and finding whether a given pair of variables might share the same memory location on entry to a given procedure.

Interprocedural optimization refers to program transformations that involve more than one procedure in a program. The most familiar example of an interprocedural optimization is *inlining*, by which the body of a procedure is substituted at a point of invocation. A broader view might consider any transformation that uses interprocedural knowledge—gathered by an interprocedural analyzer—as an interprocedural optimization. For the purposes of this paper, we will adhere to the more narrow view.

2. Interprocedural Analysis

To illustrate the need for the results of interprocedural analysis, we will introduce several important problems through a series of examples.

2.1. Interprocedural Problems

Modification and Reference Side Effects

We begin with a simple vectorization problem. Suppose we have the following code fragment:

```

COMMON X, Y
...
DO 100 I = 1, N
S0      CALL S
50      X(I) = X(I) + Y(I)
100 CONTINUE

```

Without some sort of interprocedural optimization, it is impossible to vectorize the call, but might it be possible to vectorize the assignment at statement 50? Since both X and Y are in `COMMON`, we must be concerned with side effects to these variables due to the call at S_0 . A statement can be vectorized if it is legal to distribute the loop around the two statements S_0 and 50. In turn, this is possible if there is no cycle of dependences involving both statements. We can be sure that there is no such dependence cycle if the call meets the following criteria:

1. it neither uses nor modifies X and
2. it does not modify Y .

The first condition ensures that there can be no dependence cycle involving both the call and statement 50 due to variable X . The second condition rules out a dependence cycle due to Y .

To address this problem, we introduce the interprocedural `MOD` and `REF` problems.

Definition 1. At a given call site s , the *modification side effect set* $\text{MOD}(s)$ is the set of all variables that may be modified as a side effect of the call at s . The *reference side effect set* $\text{REF}(s)$ is the set of all variables that may be referenced as a side effect of the call at s .

With these definitions, we can now formally restate the condition under which the assignment statement above can be vectorized, namely:

$$X \notin \text{REF}(S_0) \text{ and } X \notin \text{MOD}(S_0) \text{ and } Y \notin \text{MOD}(S_0).$$

Alias Analysis

Suppose we have the following subroutine:

```

SUBROUTINE S(A,X,N)
  REAL A(*), X, Y
  COMMON Y
  DO 100 I = 1, N
S0      X = X + Y*A(I)
100      CONTINUE
END

```

In compiling this loop for any machine, it would be efficient to keep both variables x and y in registers throughout the duration of the loop, without storing x until the loop is completed. This looks straightforward, but what if subroutine s is invoked by a call like `CALL S(A,Y,N)`? This creates a situation where s has two names for the variable y . We say that x and y are *aliases* in s when it is invoked in this way. Because x and y are aliased, not storing x on every iteration may overlook necessary updates to y within the loop¹.

To avoid problems like this one, the compiler must determine when two variables might be aliased to one another. In this paper, we limit our attention to aliases that arise from interactions of call-by-reference formal parameters and global variables. Because these relationships change only at call sites, we can capture these effects by computing the aliases that might hold on entry to a given subroutine.

Definition 2. For a procedure p and a formal parameter x passed to p , the *alias set* $\text{ALIAS}(p,x)$, is the set of variables that may refer to the same location as x on entry to p .

In the example above, x and y may be kept in registers without storing to memory if $y \notin \text{ALIAS}(s,x)$.

Call Graph Construction

The call graph of a program models the calling relationships between procedures in a program.

Definition 3. The *call graph* of a program is a graph $G = (N, E)$ where the vertices in N represent procedures in the program and the edges in E represent possible calls.

It is common for each distinct call site in a program to be represented by a distinct edge, in which case the call graph is actually a multigraph. We will adopt this convention in the remainder of the paper.

The accuracy of the call graph directly affects the precision of the data-flow information produced. However, construction of a precise call graph is in itself an interprocedural analysis problem. In a language in which each call must be to a named constant procedure, the call graph is easy to construct—you need only examine the body of each procedure p , entering for each call site s in p an edge (p,q) to the procedure q called at s .

However, a precise call graph is more difficult to construct in a language that permits procedure variables. Even in Fortran, where there are no assignable procedure variables, problems arise due to *procedure parameters*—formal parameters which may be bound to procedure names at the point of call. Consider the following example:

¹Knowledgeable readers will comment that the Fortran standard defines this usage to be illegal, saying that if two variables are aliases on entry to a loop then the program is not standard-conforming if the subroutine stores into either. Although this is an easy escape from the specific problem presented, it only works for Fortran, as C has no such prohibition.

```

SUBROUTINE S(X,P)
S0    CALL P(X)
      RETURN
      END

```

The question is: What procedure or procedures may be called at statement S_0 ? In other words, what values can P have at this call site? We could attempt to answer this question by examining the actual parameters at all the call sites for subroutine S , but any actual parameter could itself be a procedure parameter. Thus we must propagate procedure constants through the call graph before we can finish building it.

Definition 4. For a given procedure p and call site s within p , the *call set* $CALL(s)$ is the set of all procedures that may be invoked at s .

Although it is stated as a property of call sites, $CALL(s)$ is not a side-effect problem. It really asks what procedures can be passed to formal parameters on entry to the procedure containing s . Because it depends on the context in which the procedure containing s is invoked, it resembles alias analysis more than the side-effect problems.

Live and Use Analysis

An important data-flow problem that has been studied extensively in the literature of single-procedure analysis is the analysis of *live variables*. A variable x is said to be *live* at a given point s in a program if there is a control-flow path from s to a use of x that contains no definition of x prior to the use.

One important application of live analysis is in determining whether a private variable in a parallel loop needs to be assigned to a global variable at the end of the loop's execution. Consider the following code fragment:

```

DO I = 1, N
  T = X(I)*E
  A(I) = T + B(I)
  C(I) = T + D(I)
ENDDO

```

This loop can be parallelized by making T a local variable in the loop. However, if T is used later in the program before being redefined, the parallelized program must assign the last value of the local version of T to the global version of T . In other words, the code must be transformed as follows:

```

PARALLEL DO I = 1, N
  LOCAL t
  t = X(I)*E
  A(I) = t + B(I)
  C(I) = t + D(I)
  IF (I.EQ.N) T = t
ENDDO

```

We use the typographic convention that compiler-generated variables are represented in the code by small letters.

This code could be simplified if we could determine that variable T is not live on exit from the original loop. In this case, the conditional at the end of the loop could be eliminated to produce:

```

PARALLEL DO I = 1, N
  LOCAL t
  t = X(I)*E
  A(I) = t + B(I)
  C(I) = t + D(I)
ENDDO

```

Although live analysis can itself be viewed as an interprocedural problem, it is conveniently dealt with in terms of another interprocedural problem. *Use analysis* is the problem of determining whether there is an *upwards-exposed use* of a variable on some path through a procedure invoked at a particular call site s . An upwards exposed use is one that is not preceded by a definition of the variable on some path to the use from the point of invocation.

Definition 5. For a given call site s , which invokes procedure p , the *use side effect set* $USE(s)$ is the set of variables that have an upwards-exposed use in p .

Given this definition, we can give a more formal specification of when a variable is live. A variable x is *live* at a call site s if $x \in USE(s)$ or if there is a path through the procedure called at s that does not assign a new value to x and x is live at some control-flow successor of s .

Kill Analysis

Problems like REF, MOD and USE ask questions about what might happen on some path through a called subprocedure. It is often useful to ask about what *must* happen on every path through a procedure. The following example illustrates this:

```

L      DO I = 1, N
S0    CALL INIT(T,I)
        T = T + B(I)
        A(I) = A(I) + T
      ENDDO

```

There are two problems that might keep this loop from being correctly parallelized. First, not knowing what the subroutine INIT does, we must assume that it assigns variables in a way that creates a cycle of dependences. One way to do this would be to use, then assign, a variable that is global to the program. For example, if INIT were defined as in the following code, parallelization would be precluded:

```

SUBROUTINE INIT(T,I)
  REAL T
  INTEGER I
  COMMON X(100)
  T = X(I)
  X(I+1) = T + X(1)
END

```

Here INIT creates a cycle of dependences involving the COMMON array X, with respect to the loop in the calling program.

However, even if we can prove that the loop does not modify any global variables or even any static local (i.e., “SAVE”) variables, the call presents a more subtle problem. If this loop is to be parallelized, it must be possible to recognize that variable T can be made private to each iteration. This is possible, for example, if subroutine INIT is simply there for the purpose of initializing T, as in:

```

SUBROUTINE INIT(T,I)
  REAL T
  INTEGER I
  COMMON X(100)
  T = X(I)
END

```

Here T is initialized before being used on every iteration of the loop. Thus it is not involved in any kind of carried dependence within the loop and may be made private.

How can we discover this fact? Certainly, MOD can tell us that no global variables are modified within the subroutine and we assume that a similar analysis could be used to preclude the possibility that any static local variable in INIT is used before being modified. Therefore, the key to determining that the loop can be parallelized is to establish that variable T is assigned before being used *on every path* through INIT.

The problem of discovering whether a variable is assigned on every path through a called procedure is known as KILL, because an assignment is said to “kill” a previous value of the variable.

Definition 6. For a given call site s , the *kill side effect set* $KILL(s)$ is the set of variables that are assigned on every path through the procedure p invoked at s and any procedures invoked from within p .

Assuming that there are no global variables in $MOD(s_0)$ and that INIT does not use any static local variables before they are assigned, then the loop L can be parallelized if variable T is killed on every path through INIT and there is no path into INIT on which a use appears before any kill. This can be expressed formally as

$$T \in (\text{KILL}(S_0) \cap \neg \text{USE}(S_0)).$$

Assuming that call site s is in a block by itself, it is possible to define $\text{LIVE}(s)$ within the procedure containing s as follows:

$$(2.1) \quad \text{LIVE}(s) = \text{USE}(s) \cup (\neg \text{KILL}(s) \cap \bigcup_{b \in \text{succ}(s)} \text{LIVE}(b)).$$

where $\text{succ}(s)$ is the set of successors in the control-flow graph for the block containing s .

This generalizes the live computation to the interprocedural case if we know, by some additional analysis, the set of variables that are live on exit from the procedure. Note that the domain of $\neg \text{KILL}(s)$ is the set of variables visible in the procedure that contains s .

Constant Propagation

Constant propagation, one of the most important problems in single-procedure data-flow analysis, is also an important interprocedural problem. Consider the following program, which is abstracted from code in LINPACK:

```

SUBROUTINE S(A,B,N,IS,I1)
  REAL A(*), B(*)
L      DO I = 0, N-1
S0      A(IS*I + I1) = A(IS*I + I1) + B(I+1)
        ENDDO
      END

```

If we wish to vectorize loop L in this subroutine, a problem arises because the variable IS might take on the value 0. If that were so, there would be an output dependence. In that case, statement S_0 would actually be a reduction and could not be vectorized using the usual techniques. Although we could test for this situation at run time, we can avoid it entirely if we can determine that, on every invocation of subroutine S in the program containing it, the value of IS is always 1 (the most common case).¹

Definition 7. Given a program and a procedure p within that program, the set of *interprocedural constants* $\text{CONST}(p)$ contains the variables that have known constant values on every invocation of p . For a variable $x \in \text{CONST}(p)$, $\text{valin}(x, p)$ is a function returning the value of x on entry to p .

Although constant propagation is intractable even in a single procedure, approximate solutions in single procedures have been shown to be effective [20].

¹Barring that, it would be vectorized if it can be established that $\text{IS} \neq 0$ every time the loop is entered. Analysis of predicates like this can also be handled by a variation of constant propagation.

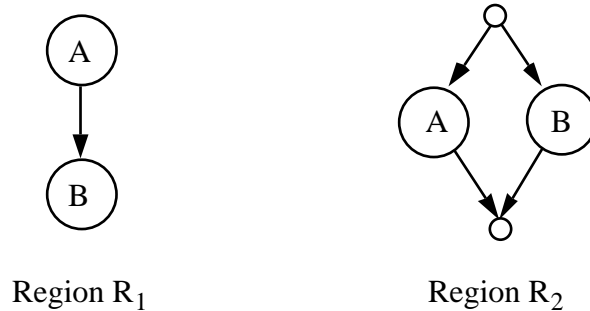


FIGURE 1. Call graph subregions.

Similarly, approximate solutions of the interprocedural constant propagation problem have been shown to determine many facts that are useful for optimization and parallelization [34, 41].

2.2. Interprocedural Problem Classification

We will now explore various classifications of interprocedural data-flow problems. These classifications are useful, because problems in the same class can generally be solved by employing similar algorithmic approaches.

May and Must Problems

There is a distinction between problems that ask whether some event “may” happen or “must” happen. MOD, REF and USE are *may problems* because they compute sets of variables that may be modified, may be referenced or may be used before being defined, respectively. On the other hand, KILL is a *must problem* because it computes a set of variables that must be killed.

Although this distinction has been extensively discussed in the literature, it is not a deep one because the converse of every may problem is a must problem and vice versa. For example, $\neg\text{MOD}(s)$ is the set of all variables that are not in MOD(s). Therefore, the $\neg\text{MOD}$ problem seeks to find those variables that must not be modified as a side effect of a given call site. Thus $\neg\text{MOD}$ is a must problem. Similarly, $\neg\text{REF}$ is also a must problem. On the other hand, $\neg\text{KILL}$ seeks to find for each call site s in the program the set of variables that may not be modified on some path through the procedure being invoked. Thus the converse of KILL, a must problem, is a may problem. Since the solution to any set problem may be converted to the solution of its converse by subtraction from the universal set, which takes time linear in the number of answer sets, there is no difference in the complexity of a must problem and its corresponding may problem.

Flow-Sensitive and Flow-Insensitive Problems

Banning introduced a seemingly related notion of flow-sensitive and flow-insensitive problems. Intuitively, a *flow-sensitive* problem is one whose precise solution requires tracing individual control-flow paths through the body of the

called subroutine. A *flow-insensitive* problem, on the other hand can be solved by examination of the body of the called subroutines without regard to control flow. Thus, MOD and REF are flow-insensitive problems because, assuming all code in a subroutine is reachable, any modification of a variable x in the body of a subroutine p means that $x \in \text{MOD}(s)$ for any call site s that invokes p . On the other hand, KILL is a flow-sensitive problem because, for a given variable x , its solution requires checking every path through the procedure body to ensure that it contains a definition of x .

Banning defines flow-sensitive and flow-insensitive problems more formally, based on how solutions on subregions of a call graph would be composed into solutions for a larger region [5] Suppose we have MOD for each subregion A and B of regions R_1 and R_2 . Then we can compose these into solutions for the whole regions as follows:

$$\text{MOD}(R_1) = \text{MOD}(A) \cup \text{MOD}(B)$$

$$\text{MOD}(R_2) = \text{MOD}(A) \cup \text{MOD}(B)$$

This is because in either decomposition, a variable may be modified if it may be modified in either subregion. There is a similar pair of equations for REF:

$$\text{REF}(R_1) = \text{REF}(A) \cup \text{REF}(B)$$

$$\text{REF}(R_2) = \text{REF}(A) \cup \text{REF}(B)$$

Now let us carry out the same exercise for KILL. In the case of region R_1 , a variable is killed if it is killed in the first region or if it is killed in the second region:

$$\text{KILL}(R_1) = \text{KILL}(A) \cup \text{KILL}(B)$$

In the case of region R_2 , a variable must be killed in both regions, to be in $\text{KILL}(R_2)$:

$$\text{KILL}(R_2) = \text{KILL}(A) \cap \text{KILL}(B)$$

The equations for USE are equally easy to develop:

$$\text{USE}(R_1) = \text{USE}(A) \cup (\neg \text{KILL}(A) \cap \text{USE}(B))$$

$$\text{USE}(R_2) = \text{USE}(A) \cup \text{USE}(B)$$

Only this second equation is the same as the corresponding equation for MOD and REF.

In examining these equations, we notice that the equations for MOD and REF use only set union as a connector, while the equations for KILL and USE are more complicated, using other connectors and often other local sets. This leads us to the following definition.

Definition 8. An interprocedural data-flow problem is *flow-insensitive* if and only if, on the parameter-free version of the problem, the value of the solution on both sequentially and alternatively composed regions (R_1 and R_2 of Figure 1) is determined by taking the union of the solutions for the subregions.

The literature on optimization sometimes refers to *flow-insensitive analysis* of a flow-sensitive problem. Our interpretation of this term is that it refers to approximating a flow-sensitive problem with a flow-insensitive one. As an example, suppose you wished to approximate a solution to the USE problem with the solution to one or more flow-insensitive problems. Note that most optimizations will be performed when it is known that a variable is not in $USE(s)$ for a particular call site s in the procedure being optimized. Therefore, we wish any approximation $APUSE$ to be conservative in the sense that it contain all of USE:

$$USE(s) \subseteq APUSE(s) \text{ or } \neg APUSE(s) \subseteq \neg USE(s)$$

If this is the case we will never depend on an untrue fact. One possible approximation is given by:

$$APUSE(s) = REF(s)$$

which is clearly a superset of $USE(s)$ and can be computed by solving a flow-insensitive problem. This might be useful because, as we shall see, flow-insensitive problems are easier to solve than flow-sensitive ones. However, we do not recommend the use of approximations like this one because more precise approximations can be computed in reasonable time using flow-sensitive analysis.

Side Effect Problems versus Propagation Problems

A final classification of interprocedural problems is by direction of data flow. Some interprocedural problems ask what may or must happen as a side effect of a procedure call that is about to be made. This class, which we refer to as *side-effect problems*, includes MOD, REF, KILL, and USE. Problems of this sort are analogous to the backward data-flow problems of single-procedure analysis.

A second class of problems asks what conditions may or must hold upon entry to the current procedure (which presumably we are interested in optimizing). We call these *propagation problems*. ALIAS and CALLS are propagation problems, as is a simple version of CONST. (A full version of CONST accounts for values returned by functions and subroutines; thus it has data-flow in both directions.) The method we developed for classifying side-effect problems as flow-sensitive or

flow-insensitive will not work for propagation problems, because these problems look back up the call chain. However, it is generally agreed that the ALIAS problem for Fortran is flow-insensitive. Control flow inside the procedure is unimportant, because the set of aliases introduced at a call site is a local, textual property of the call site itself. Events that happen between procedure entry and the call site cannot change the set of aliases introduced by the call.

On the other hand, the constant propagation problem for Fortran is flow-sensitive because, for a constant to propagate to a called procedure, every path through the calling procedure (and its preceding interprocedural context) must result in the same constant value arriving at the call site. By the same argument, alias analysis for languages like C with pointer assignment is also flow-sensitive.

A Problem Classification Table

The breakdown of problems in the two dimensions is summarized in Table 1 below.

	propagation	side effect
flow-insensitive	ALIAS, CONST	MOD, REF
flow-sensitive	CONST	KILL, USE

TABLE 1. Interprocedural problem classification

The classification into flow-sensitive and flow-insensitive is an important one because Banning established, by providing an algorithm, that flow-insensitive problems can be solved in time that is polynomial in the size of the call graph [5], while Myers showed that flow-sensitive problems are intractable in the presence of aliasing, recursion, and arbitrary nesting [57].

2.3. Flow-Insensitive Side-Effect Analysis

We now turn to the solution of the flow-insensitive side-effect analysis problems. Throughout this section, we use the modification side-effect analysis (MOD) problem as our example. The reference side-effect analysis (REF) problem can be solved using exactly the same approach.

Assumptions

We begin by establishing a set of assumptions that are representative characteristics of Fortran and, in some cases, C. First, we assume that there is no procedure nesting—that is, variables are subdivided into local and global sets. Local variables are known only within their home procedure and global variables are known in every procedure. Although this might seem restrictive, the methods we will introduce here have easy extensions to the case of general nesting.

A second assumption is that all parameters are passed by reference and that there are no pointer variables. The purpose of this restriction is to simplify the aliasing patterns that the algorithms encounter. Under this restriction, aliases

can only be introduced at call sites. Although both Fortran 90 and C have pointer variables, this assumption is valid for Fortran 77, which is the input language for most automatic parallelization systems.

Even though Fortran 77 does not support recursion, the algorithms presented here work correctly in the presence of recursive procedure calls. We do however make an assumption about the size of parameter lists to procedures, namely that the maximum number of formal parameters to a procedure does not grow with the size of the program. In other words, programmers do not typically deal with the increasing complexity of a large program by increasing the complexity of the procedure interfaces. Thus we will assume that there exists a constant μ , such that the number of formal parameters in any procedure p is less than or equal to μ .

MOD Problem Formulation

The goal is to compute, for each call site s in the program, the set $\text{MOD}(s)$ that contains every variable that may be modified as a side effect of the procedure call at s . Our first step is to note that we can simplify the problem by disregarding potential aliases. Specifically, we will compute the set $\text{DMOD}(s)$, the *direct modification side-effect set*, which contains all variables visible at s that are directly modified as a side effect of the call. $\text{DMOD}(s)$ may be smaller than $\text{MOD}(s)$ because it does not take into account the fact that a variable x that is directly modified as a side effect may have several possible aliases at the point of call. Each of these aliases must be in the final $\text{MOD}(s)$ if the solution is to be precise. However, once $\text{DMOD}(s)$ has been computed it can be updated to $\text{MOD}(s)$ with the help of the ALIAS sets. Recall that for a given procedure p , $\text{ALIAS}(p, x)$ contains the set of all variables that may be aliased to x on entry to p . Given these sets, we can update $\text{DMOD}(s)$ to $\text{MOD}(s)$ according to the following formula:

$$(2.2) \quad \text{MOD}(s) = \text{DMOD}(s) \cup \bigcup_{x \in \text{DMOD}(s)} \text{ALIAS}(p, x)$$

where p is the procedure containing call site s . Construction of $\text{ALIAS}(p, x)$ will be discussed in Section 2.4.

The computation of $\text{DMOD}(s)$ is formulated in terms of another set, the *generalized modification side effect set*, $\text{GMOD}(p)$. The difference between these sets is conceptually simple. The domain of $\text{DMOD}(s)$ is the name space of the calling procedure, while the domain of $\text{GMOD}(p)$ is the name space of the called procedure. Thus, $\text{GMOD}(p)$ contains the set of global variables and formal parameters of p that might be modified if p is invoked. This modification can occur directly, inside p , or indirectly, as a result of some call made by p . Given $\text{GMOD}(p)$ for each procedure in the program, we can compute $\text{DMOD}(s)$ by the following formula:

$$(2.3) \quad \text{DMOD}(s) = \{z \mid s \text{ invokes } p, z \xrightarrow{s} w \text{ and } w \in \text{GMOD}(p)\},$$

where $z \xrightarrow{s} w$ means that, at call site s , actual parameter z is bound to formal parameter w of the called procedure. Equation 2.3 reduces the problem to computing, for every procedure p in the program, the set $\text{GMOD}(p)$ of variables that may be changed as a side effect of invoking p . Generally, $\text{GMOD}(p)$ will contain two types of variables:

1. those that are explicitly modified in the body of procedure p and
2. those that are modified as a side effect of some procedure that p invokes.

Let $\text{IMOD}(p)$, the *immediate modification side-effect set*, denote the set of variables explicitly modified in p . Then the following formula holds:

$$(2.4) \quad \text{GMOD}(p) = \text{IMOD}(p) \cup \bigcup_{s=(p,q)} \{z | z \xrightarrow{s} w \text{ and } w \in \text{GMOD}(q)\}$$

Note that the union is taken over all call sites s within p .

Problem Decomposition

The system of data-flow equations in Equation 2.4 can be solved by the iterative method of data-flow analysis. However, the solution may take a long time to converge because the system does not satisfy the conditions under which the iterative method can be guaranteed to converge rapidly [52]. Furthermore it is not formulated in a way that permits the use of a fast elimination method [33]. Fast solution methods rely on being able to limit the number of times that the analysis must traverse a loop in the problem graph. For flow-insensitive interprocedural analysis, the problem graph is the call graph, so convergence problems are limited to recursive regions.

A closer examination of the problem with recursion reveals that it is related to reference formal parameters. Consider the following example:

```

SUBROUTINE P(F0,F1,F2,...,Fn)
  INTEGER X,F0,F1,F2,...,Fn
  ...
S0      F0 = some expression
  ...
S1      CALL P(F1,F2,...,Fn,X)
  ...
END
```

In this example it is possible to see why, in the general case, the analyzer must iterate around the recursive cycle an unbounded number of times. The question being asked is: How many of the parameters of subroutine P may be modified as a side effect of invoking P ? Clearly, F_0 can be modified at statement S_0 , but we have to examine the recursive call at S_1 to discover that F_1 is passed to F_0 , so it too can be modified. One more time around the recursive loop reveals

that F_2 can also be modified. This procedure continues until it is discovered that F_n —the last parameter—may also be modified. If n is unbounded, the number of iterations over the recursive cycle is unbounded as well.

These observations make it clear why we assumed an upper bound on the number of parameters to any subroutine—it permits us to establish a constant upper bound on the number of iterations required for the process to converge. However, we can achieve a better time bound by decomposing the problem further, treating side effects to reference parameters separately from side effects to global variables. This is achieved by introducing an extended version of the immediate modification side effect set $\text{IMOD}(p)$, called $\text{IMOD}^+(p)$, which will contain all of $\text{IMOD}(p)$ plus all those variables that may be modified as a result of side effects to reference formal parameters of procedures invoked from within p . In other words, a variable x is in $\text{IMOD}^+(p)$ if

1. $x \in \text{IMOD}(p)$ or
2. $x \xrightarrow{s} z$ and $z \in \text{GMOD}(q)$, where $s = (p, q)$.

If we can compute $\text{IMOD}^+(p)$ for every procedure p in the program, we can solve for $\text{GMOD}(p)$ using the following simple system of equations:

$$(2.5) \quad \text{GMOD}(p) = \text{IMOD}^+(p) \cup \bigcup_{s=(p,q)} \text{GMOD}(q) \cap \neg \text{LOCAL}$$

where LOCAL is the set of all local variables in the program, so its complement is the set of all global variables. Since all side effects to reference formal parameters of p are reflected in $\text{IMOD}^+(p)$, we need only be sure that all side effects to global variables are added by the union of the GMOD sets for successors. If a global variable is modified as a side effect of invoking p it must be modified directly in the text of the subroutine, in which case it is in $\text{IMOD}(p) \subseteq \text{IMOD}^+(p)$, or it is passed as an actual parameter to another subroutine where it is modified, in which case it is in $\text{IMOD}^+(p)$ by definition, or it is modified as a global by some subroutine called directly or indirectly from p , in which case it must be in $\text{GMOD}(q)$ for some successor q of p , which establishes Equation 2.5.

Thus, we have decomposed the problem into two parts:

1. the computation of $\text{IMOD}^+(p)$ for every procedure p in the program and
2. the propagation of global information according to Equation 2.5.

The next two subsections describe these calculations in more detail.

Solving for IMOD^+

To compute the IMOD^+ sets for each procedure p , we will formulate a final set, called $\text{RMOD}(p)$ that contains the formal parameters in procedure p that may be modified in p , either directly or as a side effect of some call that p makes. Given

$\text{RMOD}(p)$ for each procedure p in the program, the IMOD^+ can be computed using an equation that is a direct analog of Equation 2.4:

$$(2.6) \quad \text{IMOD}^+(p) = \text{IMOD}(p) \cup \bigcup_{s=(p,q)} \{z | z \xrightarrow{s} w \text{ and } w \in \text{RMOD}(q)\}$$

Thus, we can construct $\text{IMOD}^+(p)$ by first constructing $\text{RMOD}(p)$ for each procedure and then applying Equation 2.6.

To construct the RMOD sets, we introduce a new data structure, which we call the *binding graph*. It directly represents the parameter-to-parameter bindings that can occur in the program. It is built as follows:

1. for every formal parameter f of every procedure in the program, create a vertex in the binding graph
2. if f_1 is a formal parameter of procedure p , f_2 is a formal parameter of procedure q , and there is a call site $s = (p, q)$ where $f_1 \xrightarrow{s} f_2$, then add an edge (f_1, f_2) to the binding graph.

An important question to consider is: How large is the binding graph? If N is the number of vertices in the call graph and E is the number of edges, then the number of vertices in the binding graph can be no more than μN , where μ is the upper bound on the number of parameters to any procedure in the program. This is true because there can be no more vertices than the number of formal parameters in the entire program, which is clearly bounded by μN . Similarly, each call graph edge can give rise to no more than μ edges, since no more than one formal reference parameter can appear in a single actual parameter position. Thus the total number of edges in the binding graph is no larger than μE and the binding graph is no more than a constant factor greater in size than the call graph; that is, its size is $O(N + E)$.

To construct $\text{RMOD}(p)$, we will use a simple marking algorithm on the binding graph, in which each vertex is annotated with a logical mark implementable with a single bit. Initially, all the marks are set to false. Next, for each formal parameter f of procedure p , if $f \in \text{IMOD}(p)$, the corresponding node's mark is set to true. True bits are propagated around the graph using the rule that any formal f_1 that is bound to formal f_2 with a true mark must also be marked true. When no more propagation is possible, $\text{RMOD}(p)$ is the set of formals of p that are marked true. The algorithm is given in Figure 2.

Theorem 2. Algorithm *computeRmod* in Figure 2 correctly computes the RMOD sets—that is, on exit $f \in \text{RMOD}[proc[f]]$ if and only if f may be modified as a side effect of invoking the procedure $p = proc[f]$.

Define
 $proc[f]$
somewhere

Proof. If. Assume that f may be modified. Then either it is modified in $proc[f]$, in which case it is marked true in statement S_1 , or there is a path in the

```

procedure computeRmod( $P, N_B, E_B, \text{IMOD}, \text{proc}, \text{RMOD}$ )
  /*  $P$  is the collection of procedures in the program */
  /*  $N_B$  is the collection of formal parameters in the binding graph */
  /*  $E_B$  is the collection of edges in the binding graph */
  /*  $\text{mark}[f]$  maps formal parameters to their mark values */
  /*  $\text{proc}[f]$  maps a parameter to its procedure */
  /*  $\text{IMOD}[p]$  maps a procedure to its immediate MOD side effect set */
  /*  $\text{RMOD}[p]$  is the collection of output sets */
  /*  $\text{worklist}$  is a working set of formal parameters */

  L1   for each  $f \in N_B$  do  $\text{mark}[f] := \text{false}$  od;
         $\text{worklist} := \emptyset$ ;

  L2   for each  $f \in N_B$  such that  $f \in \text{IMOD}[\text{proc}[f]]$  do
  S1      $\text{mark}[f] := \text{true}$ ;
         $\text{worklist} := \text{worklist} \cup \{f\}$ 
        od;
  L3   while  $\text{worklist} \neq \emptyset$  do
         $f :=$  an arbitrary element in  $\text{worklist}$ ;
         $\text{worklist} := \text{worklist} - \{f\}$ ;
  L4     for each  $v$  such that  $(v, f) \in E_B$  do
  S2       if  $\text{mark}[v] = \text{false}$  then
  S3          $\text{mark}[v] := \text{true}$ ;
         $\text{worklist} := \text{worklist} \cup \{v\}$ 
        fi
        od
        od;
  L5   for each  $p \in P$  do  $\text{RMOD}[p] := \emptyset$  od;
  L6   for each  $f \in N_B$  do
        if  $\text{mark}[f]$  then  $\text{RMOD}[\text{proc}[f]] := \text{RMOD}[\text{proc}[f]] \cup \{f\}$  fi
        od
end computeRmod

```

FIGURE 2. Algorithm for constructing RMOD sets.

binding graph to a formal parameter f_0 which is in $\text{IMOD}[\text{proc}[f_0]]$. Thus f_0 is marked true and added to the worklist in loop L₂. Let $f = f_n, f_{n-1}, \dots, f_1, f_0$ be the sequence of parameters in the binding graph that constitute the path from f to f_0 . Suppose $f = f_n$ is never marked true. Then there must be a minimum k such that f_k is never marked true by the algorithm but f_{k-1} is marked true. But since f_{k-1} is put on the worklist when it is marked true and it is taken off the worklist eventually and every incoming edge to f_{k-1} is examined at that time, f_k must be marked true, a contradiction. Thus every parameter that may be modified as a side effect of invoking $\text{proc}[f]$ is put into $\text{RMOD}[\text{proc}[f]]$.

Only if. Suppose f is a parameter that is marked true but cannot be modified. The algorithm only marks formal parameters true if they are modified within the body of their procedure or if there is a path in the binding graph to another parameter which is modified in its procedure. Since the binding graph has an edge only if the corresponding procedure call exists, it must be the case that the original parameter f may be modified. \square

Theorem 3. In the worst case, algorithm *computeRmod* requires $O(N + E)$ steps, where N and E are the number of vertices and edges in the call graph, respectively.

Proof. Let N_B and E_B be the number of vertices and edges, respectively in the binding graph. We know, by definition, that $N_B \leq \mu N$ and $E_B \leq \mu E$. Thus, if we can show that the algorithm runs in time proportional to the size of the binding graph, we will have established the result. Clearly loop L_1 requires $O(N_B)$ steps. If the membership test on IMOD takes constant time and adding an element to the worklist takes constant time, then loop L_2 also runs in $O(N_B)$ steps. Loop L_5 takes time proportional to the number of procedures in the program times the time to initialize the RMOD sets. If these sets are implemented as bit vectors of length μ then this takes $O(N \cdot \log_2 \mu)$ time. Loop L_6 takes $O(N_B)$ assuming that adding to RMOD takes constant time. (As an existence proof, we can meet the time constraints by using bit-vector implementations for IMOD and RMOD and a linked list for the worklist.)

All that remains is to analyze the running time of loops L_3 and L_4 . If we assume that choosing an arbitrary element from the worklist takes constant time, then L_3 makes N_B iterations, since a vertex appears on the worklist at most once. If the edges are arranged as predecessor lists, then loop L_4 requires time proportional to the number of predecessors. Taken over all the vertices in the graph, this takes at most $O(E_B)$ time.

Since $N_B \leq \mu N = O(N)$ and $E_B \leq \mu E = O(E)$, we have established that the algorithm takes $O(N + E)$ steps. \square

Given RMOD sets for each procedure, we can construct $\text{IMOD}^+(p)$ by visiting each call site $s = (p, q)$ and each parameter at s to determine if it is bound to a variable in $\text{RMOD}(q)$. To construct IMOD^+ sets, we apply Equation 2.6. The N IMOD^+ sets must be initialized; they have at most $V + \mu$ elements; this entire process takes $O(NV)$ time. If the membership test in RMOD is $O(1)$, adding in the elements contributed by the various RMOD sets takes $O(\mu E) = O(E)$ time. Thus, it takes $O(NV + E)$ time to compute the IMOD^+ sets, including initializations.

Solving for GMOD

Once we have constructed $\text{IMOD}^+(p)$ for each procedure p in the program, we must use it to compute $\text{GMOD}(p)$ according to Equation 2.5, which we repeat here:

```

procedure findGmod(N, E, n,  $\text{IMOD}^+$ , LOCAL)

  integer dfn[n], lowlink[n], nextdfn, p, q, d,
     $\text{IMOD}^+$ [n], GMOD[n], LOCAL;
  integer stack Stack;

  procedure search(p);
    dfn[p] := nextdfn; nextdfn := nextdfn + 1;
    GMOD[p] :=  $\text{IMOD}^+$ [p]; lowlink[p] := dfn[p];
    push p onto Stack;

    for each q adjacent to p do
      if dfn[q] = 0 then /* tree edge */
        search(q);
        lowlink[p] := min(lowlink[p], lowlink[q])
      fi;

      if dfn[q] < dfn[p] and q ∈ Stack then
        lowlink[p] := min(dfn[q], lowlink[q])
      else /* apply equation */
        GMOD[p] := GMOD[p] ∪ (GMOD[q] ∩ ¬LOCAL)
      fi
    od;

    /* test for root of strong component */
    if lowlink[p] = dfn[p] then

      /* adjust GMOD sets for each member of the SCR */
      repeat
        pop u from Stack;
        GMOD[u] := GMOD[u] ∪ (GMOD[p] ∩ ¬LOCAL)
      until u = p

      fi
    end search;

    /* subroutine body */
    nextdfn := 1; dfn[*] := 0; Stack := ∅;
    search(1); /* by convention root = 1 */
  end findGmod

```

FIGURE 3. Algorithm for propagation of global modification side effects.

```

L1 for each call site  $s$  do
S1   MOD[ $s$ ] := DMOD[ $s$ ];
L2   for each  $x \in \text{DMOD}[s]$  do
L3     for each  $v \in \text{ALIAS}[\text{proc}[s], x]$  do
S2       MOD[ $s$ ] := MOD[ $s$ ]  $\cup v$ 
      od
    od
  od

```

FIGURE 4. Naive alias update of DMOD to produce MOD.

$$\text{GMOD}(p) = \text{IMOD}^+(p) \cup \bigcup_{s=(p,q)} \text{GMOD}(q) \cap \neg \text{LOCAL}$$

This equation implies that a variable x is in $\text{GMOD}(p)$ if it is in $\text{IMOD}^+(p)$ or x is global and there is a nonempty path in the control-flow graph from p to procedure q where $x \in \text{IMOD}^+(q)$. In other words we have reduced the problem to a variant of the reachability problem in the call graph. It is well known that reachability can be solved in time linear in the size of the problem graph using a depth-first search algorithm based on Tarjan's algorithm for finding strongly-connected components of a directed graph. Algorithm *findGmod* in Figure 3 accomplishes this.

Since this is a direct adaptation of depth-first search, it runs in $O(N + E)$ steps, where each step may involve a bit-vector operation of length V , the number of variables in the program. Thus the algorithm takes $O((N + E)V)$ elementary steps in the worst case.

Theorem 4. Algorithm *findGmod* in Figure 3 correctly computes the $\text{GMOD}(p)$ for every procedure p in the program.

Proof. This is a direct adaptation of Tarjan's algorithm for finding strongly-connected components. However, as the algorithm backs up in reverse invocation order, instead of collecting a set of strongly connected regions, it updates the $\text{GMOD}(p)$ set at each node. When it reaches the head of a strongly-connected component, it updates $\text{GMOD}(u)$, for every u in the component, to include the nonlocal part of the GMOD set for the head. Thus, the global parts of $\text{GMOD}(u)$ for every procedure u in the strongly connected region are identical, as they should be. The GMOD sets for nodes that are not in loops are correct by virtue of the order of visit and the GMOD sets for nodes that are in loops are correct because of the updates of all procedures in the cycle. \square

Taken together the results from this section and the previous one establish that the entire computation can be done in $O((N + E)V)$ steps. Since DMOD can be computed from GMOD in $O(NV + E)$ time, the complete computation of DMOD also takes $O((N + E)V)$ steps. Asymptotically, this cannot be improved,

for each call site s in the program **do**

```

     $t := \emptyset;$       /*  $t$  is a temporary bit vector of length  $\mu^*$  */
    for each global variable  $x \in \text{DMOD}[s]$  do    /*  $O(V)$  iterations */
         $t := t \cup \{x\};$                         /* constant time */
         $t := t \cup \text{ALIAS}[\text{proc}[s], x]$         /*  $O(\mu)$  time */
    od;

```

```

     $\text{MOD}[s] := \text{MOD}[s] \cup t;$                     /*  $O(V)$  time */

```

```

    for each formal parameter  $f \in \text{DMOD}[s]$  do /*  $O(\mu)$  iterations */
         $\text{MOD}[s] := \text{MOD}[s] \cup \{f\};$             /* constant time */
         $\text{MOD}[s] := \text{MOD}[s] \cup \text{ALIAS}[\text{proc}[s], f]$  /*  $O(V)$  time */
    od

```

od

FIGURE 5. Fast combination of side-effect and alias information.

because the analyzer must evaluate Equation 2.5 at least once at every node in the call graph, which would require $O((N + E)V)$ time.

2.4. Flow-Insensitive Alias Analysis

Update of DMOD to MOD

Once the direct modification side effect sets have been computed, there still remains the problem of factoring in aliasing. Recall the formula from Equation 2.2, which we repeat here:

$$\text{MOD}(s) = \text{DMOD}(s) \cup \bigcup_{x \in \text{DMOD}(s)} \text{ALIAS}(p, x)$$

Figure 4 shows a naive algorithm for implementing this conversion. Loop L_1 executes $O(E)$ times. The cost of each execution of S_1 is proportional to $|\text{MOD}(s)|$. The MOD set for a call site can contain the names of every global variable and all the formal parameters to the calling procedure, so $|\text{MOD}(s)|$ is $O(V)$. Loop L_2 iterates $O(V)$ times per iteration of L_1 . The body of L_2 does one set operation for each alias involving an element of $\text{DMOD}(s)$; in the worst case, this entails $O(V)$ set operations. Thus, the algorithm in Figure 4 requires $O(E \cdot (V + V \cdot V))$ time, which simplifies to $O(EV^2)$ in the worst case.

If we cannot improve the running time of alias analysis, it may dominate the running time of the entire side effect analysis, making it impractical for use on programs with nontrivial aliasing patterns. However, we can achieve a significant improvement by carefully classifying aliasing patterns that may arise.

First, we note that, in a two-level naming hierarchy, containing only global and local variables, two global variables can *never* be aliases of one another. They can only be aliased to reference formal parameters. Thus, for a global variable x , $\text{ALIAS}(p, x)$ can only contain reference formal parameters of procedure p . This means that $\text{ALIAS}(p, x)$ can be no larger than μ entries, where μ is the maximum number of formal parameters of any procedure in the program.

On the other hand, for a given formal parameter f , $\text{ALIAS}[\text{proc}[p], f]$ may contain any global or any other formal parameter. Thus, it may be of size $O(V)$.

These observations suggest that we break down the update of DMOD to MOD into two cases: one for formal parameters and the other for global variables. When considering the aliases of a particular variable $x \in \text{DMOD}(s)$,

1. if x is a global variable, we will add a very small set ($\leq \mu$ elements) to $\text{MOD}(s)$; since there are $O(V)$ global variables, we may need to do this $O(V)$ times, and
2. if x is a formal parameter of the procedure p containing s , we will add up to $O(V)$ elements to $\text{MOD}(s)$; but since there can be at most μ formal parameters, we need only do this μ times.

In each case, the amount of work is $O(V)$, not $O(V^2)$. The entire algorithm is shown in Figure 5.

It should be clear from the previous discussion that the overall running time of the combining strategy in Figure 5 is $O(V)$ per call site for a total of $O(EV)$ time. This is promising because, if we can compute the ALIAS sets in $O((N + E)V)$, the time bound for the entire MOD computation will be $O((N + E)V)$.

Computing Aliases

We now turn to the problem of computing the sets $\text{ALIAS}(p, x)$ for each procedure p and each variable x that is either global or a parameter of some procedure. To do this as rapidly as possible, we will once again take advantage of the observation that globals can only be aliases of formal parameters. The first step is to compute, for each formal parameter f in the program, an intermediate set $A(f)$. $A(f)$ is the subset of $\text{ALIAS}[\text{proc}[f], f]$ that arises from the binding of global variables to f , either directly at a site that calls $p = \text{proc}[f]$, or through some chain of bindings:

$$g \rightarrow f_0 \rightarrow f_1 \rightarrow f_2 \rightarrow \dots \rightarrow f_{n-1} \rightarrow f_n = f$$

Thus, $A(f)$ is an approximation to $\text{ALIAS}(\text{proc}[f], f)$ that contains all the global variables in that set. To compute $A(f)$, we will use a forward propagation on a variation of the binding graph in which cycles have been reduced to single nodes. The algorithm for computing $A(f)$ for each formal parameter in the program is given in Figure 6. In this algorithm, the set $A(f)$ will be represented by a bit vector of length $O(V)$.

The first phase of the algorithm performs a forward propagation of global variables over the binding graph. It requires $O((N + E)V)$ steps. The second

```

procedure computeA( $N, E, N_B, E_B, A$ )

  for each formal parameter  $f \in N_B$  do  $A[f] := \emptyset$  od;
  for each call site  $s \in E$  do
    for each global  $g$  mapped to formal parameter  $f$  at  $s$  do
       $A[f] := A[f] \cup \{g\}$ 
    od
  od;

  replace every cycle in the binding graph with a single node,
  reducing the graph to a directed, acyclic form;
  for each  $f$  in the reduced graph in topological order do

     $A[f] := A[f] \cup \bigcup_{(f_0, f) \in E_B} A[f_0]$ 
  od;
  for each cycle in the original binding graph do
    let  $C$  be the reduced binding graph node for the cycle;
    for each  $f \in C$  do  $A[f] := A[C]$  od
  od
end computeA

```

FIGURE 6. Computing approximate alias sets for formal parameters.

phase derives an acyclic graph from the binding graph and propagates the A sets along the derived graph. This takes $O(N + E)$ time to derive the graph and $O(E)$ unions on sets of size V in the propagation. The final phase maps the A set computed for the node representing a cycle back onto the individual nodes in the cycle. This requires N operations on sets of length V . Thus, the overall complexity is $O((N + E)V)$.

Once the sets $A(f)$ are available for every formal parameter in the program, the alias sets for every global variable can be computed by a simple inversion, as shown in Figure 7. This computation can be subdivided into two components: initializing the alias sets and computing the inversion. Recall that we can represent an alias set for a global variable within a given procedure as a bit vector of length μ . Therefore the initializations in step S_0 of Figure 7 take no more than $O(NV\mu) = O(NV)$ time. However, since an update can be done in constant time, the total cost of the updates in statement S_1 is also $O(NV)$.

All that remains is the computation of aliases for formal parameters. Note that $A(f)$ is an approximation for $\text{ALIAS}(p, f)$, containing all the global aliases of f . To expand $A(f)$ to $\text{ALIAS}(p, f)$, we need only add the formal parameters that may be aliased to f . In the simple two-level language we are considering, the only formals that may be aliased to f are other parameters of the same procedure. Thus there can be no more than $\mu(\mu - 1)/2$ formal pairs in any given procedure.

```

for each procedure  $p$  do
  for each global variable  $g$  do
S0      ALIAS[ $p, g$ ] :=  $\emptyset$ 
  od
od

for each formal parameter  $f$  in the program do
  for each  $g \in A[f]$  do
S1      ALIAS[ $proc[f], g$ ] := ALIAS[ $proc[f], g$ ]  $\cup \{f\}$ 
  od
od

```

FIGURE 7. Inversion to compute ALIAS[p, g] for each global variable g .

To compute the formals that may be aliased to one another, the algorithm *computePairs* in Figure 8 keeps track of the set FPAIRS(p) of pairs of formal parameters that may be aliased to one another in a given procedure p . After initializing all the FPAIRS sets to the empty set, we examine each call site for alias introductions, occurring when the same variable is passed to two different parameters. Whenever one is found, we add the formal parameter pair in the

```

procedure computePairs
   $W := \emptyset$ ;
L1   for each procedure  $p$  in the program do FPAIRS[ $p$ ] :=  $\emptyset$  od;
L2   for each alias introduction site (e.g., “CALL P(X,X)”) do
      insert the resulting formal pair in FPAIRS[ $p$ ] and
      in the worklist  $W$ ;
      od;
L3   while  $W \neq \emptyset$  do
      remove  $\langle f_1, f_2 \rangle$  from  $W$ ;
L4   for each call site  $s$  in  $proc[f_1]$  passing both  $f_1$  and  $f_2$  do
      let  $q$  be the procedure invoked at  $s$ ;
      if  $f_1$  and  $f_2$  are passed to  $f_3$  and  $f_4$  at  $s$  and
       $\langle f_3, f_4 \rangle \notin \text{FPAIRS}[q]$  then
        FPAIRS[ $q$ ] := FPAIRS[ $q$ ]  $\cup \{\langle f_3, f_4 \rangle\}$ ;
         $W := W \cup \{\langle f_3, f_4 \rangle\}$ ;
      fi
      od
      od
end computePairs

```

FIGURE 8. Algorithm to compute formal parameter pairs that may be aliased.

alias to a worklist. Then we iterate over the worklist looking for possible alias propagations, which occur when two parameters that may be aliased are both passed to another procedure. If the resulting pair of formals in the called procedure is not already in the FPAIRS set, the pair is added to the worklist. This procedure continues until the worklist is empty.

Let us now analyze the running time of the algorithm *computePairs*. Note that there can be no more than $\mu(\mu - 1)/2$ parameter pairs in any subroutine. The initialization loop L_1 takes $O(N)$ time, since it is constant time for each procedure. Since it must look at each call site, the loop at L_2 takes $O(E)$ time, but the operations are all constant-time, assuming we use some sort of linked structure for both W and FPAIRS. Loop L_3 is entered once for each pair that is put on the worklist. Since there can be at most $\mu(\mu - 1)/2$ pairs for each call site, the total number of times L_3 is entered is less than or equal to $E \cdot \mu(\mu - 1)/2$, so L_3 is entered $O(E)$ times. Similarly, each call site in the program is examined no more than $\mu(\mu - 1)/2$ times, so loop L_4 is entered $O(E)$ times. Since the body of L_4 takes constant time, the entire process takes $O(N + E)$ time.

Is this right?

Thus, we have established that alias analysis can be done in $O((N + E)V)$ time, which means that the entire MOD solution can be completed in this time.

2.5. Constant Propagation

Constant propagation is the process of discovering, at compile time, expressions that have known constant values. It has been shown to be an important and effective optimization in single procedures [20] and across whole programs [34, 41]. Unfortunately the problem is difficult even in a single procedure—obtaining a precise solution has been shown to be unsolvable [53]. Even the usual single-procedure approximate problem is flow-sensitive in the interprocedural setting, hence intractable [57].

One important reason for the difficulty is that constants propagated into a program region can make it possible to evaluate program expressions in that region, yielding new constants on exit. An example illustrating the interprocedural case is shown below.

```

SUBROUTINE PHASE(N)
  INTEGER N,A,B
  CALL INIT(A,B,N)
  CALL PROCESS (A,B)
END

SUBROUTINE INIT(A,B,N)
  INTEGER A,B,N
  A = N+1
  B = (N*A)/2
END

```

The purpose of subroutine INIT is to initialize the variables A and B. Thus if N is a constant 10 on entry to procedure PHASE, A will be a constant 11 and B will

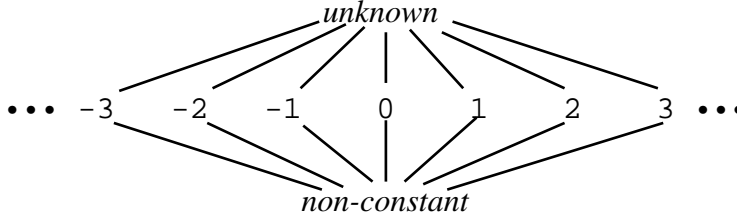


FIGURE 9. Constant propagation lattice.

be a constant 55 on exit from INIT. Hence these constant values will be available on entry to PROCESS. This section will develop an approach to determine these facts.

Figure 10 shows an algorithm for constant propagation within a single procedure, modelled after Wegman and Zadeck [72]. This algorithm uses an iterative process on the use-definition graph or the static single assignment (SSA) form of the given procedure [31].

The iterative method is guaranteed to converge because it is based on the constant propagation lattice shown in Figure 9. The algorithm in Figure 10 is linear in the size of the base graph because an instruction is only put on the worklist when its output value is lowered, which can happen only twice for each instruction in the lattice of Figure 9. In addition, each edge out of an instruction is visited once each time it is removed from the worklist, so the running time required by the algorithm is proportional to the number of nodes and edges in the graph.

To use an interprocedural analog of this iterative process, we will need to construct an interprocedural analog of the use-definition graph, which we call the *interprocedural value propagation graph*. In this graph, the vertices represent “jump functions” that compute values the values of variables passed out of a procedure (either at a call site or a return) as a function of values passed into the procedure and known constant values generated inside the procedure. Values can be propagated from the calling procedure to the called procedure; this is modelled with a “forward” jump function. Values can also be propagated when a procedure returns; this is modelled with a “return” jump function.

Let s be a call site within procedure p and x a formal parameter of the procedure q called at s . The *jump function* for x at s , denoted J_s^x , determines the value of x in terms of the values known at s . These include the values of parameters and global variables passed into p , known constant values created p , and values returned from other procedures that p calls before reaching s . Because this set of possible inputs is large, we refer to the subset that is actually used in evaluating J_s^x as the *support* of J_s^x .

We now return to the construction of the *interprocedural value graph*. Here are the steps:

1. construct a node for each forward jump function J_s^x

```

/* valin(w, s) is the current approximate value of input w to s */
/* valout(v, s) is the current value of output v from s */
/* M(s)(inputs to s) is the result of symbolic interpretation of */
/* statement s over the lattice values of its inputs. The */
/* output is the lattice value of the output of the statement. */
for all statements s in the program do
  for each output v of s do valout(v, s) := unknown od;
  for each input w of s do
    if w is a variable then valin(w, s) := unknown
    else valin(w, s) := the constant value of w fi
  od
od;
worklist := { all statements of constant form, e.g., X = 5 } ;
while worklist ≠ ∅ do
  pick an arbitrary statement x from worklist;
  let v denote the output variable for x;

  /* Symbolic interpretation of the statement x */
  newval := M(x)(valin(u, x), for all inputs u to x);
  if newval ≠ valout(v, x) then
    valout(v, x) := newval;
    for all (x, y) ∈ defuse do
      oldval := valin(v, y);
      valin(v, y) := oldval ∧ valout(v, x);
      if valin(v, y) ≠ oldval then worklist := worklist ∪ {y} fi
    od
  fi
od

```

FIGURE 10. Constant Propagation Algorithm

2. if $x \in \text{support}(J_t^y)$, where t is a call site in the procedure called at s , then construct an edge from J_s^x to J_t^y .

The iterative algorithm shown in Figure 10 can be applied to the resulting graph.

We now present a simple example of this process. Consider the program shown in Figure 11. From this example we can easily derive the jump functions:

$$J_\alpha^X = \{1\}; J_\alpha^Y = 2$$

$$J_\beta^U = \{x+y\}; J_\beta^V = \{x-y\}$$

The call graph and the resulting interprocedural value propagation graph are shown in Figure 12.

```
PROGRAM MAIN
  INTEGER A,B
  A = 1
  B = 2
 $\alpha$   CALL S(A,B)
END

SUBROUTINE S(X,Y)
  INTEGER X,Y,Z,W
  Z = X + Y
  W = X - Y
 $\beta$   CALL T(Z,W)
END

SUBROUTINE T(U,V)
  PRINT U,V
END
```

FIGURE 11. Interprocedural constant propagation example.

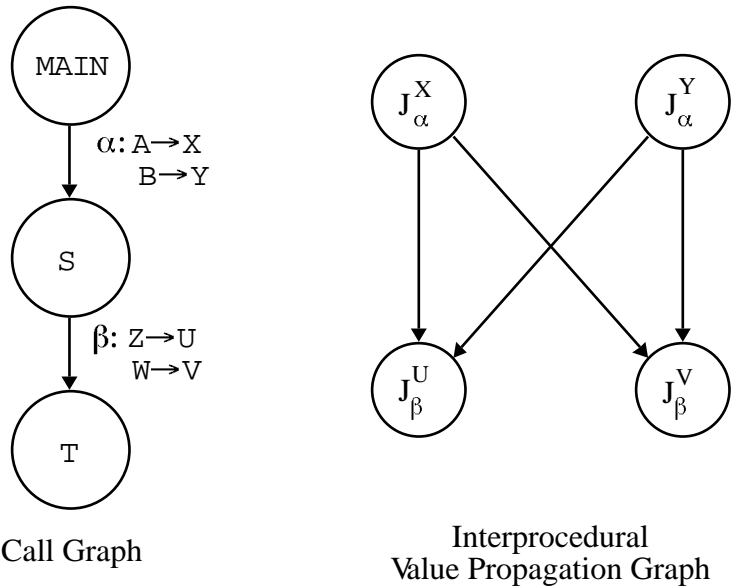


FIGURE 12. Interprocedural value propagation graph example.

It can easily be seen that the constant propagation algorithm applied to the interprocedural value propagation graph in Figure 12 will quickly converge to the constant assignments:

$$x=1; y=2; u=3; v=-1$$

To estimate the total cost of the algorithm, recall that the cost of the iterative constant propagation algorithm on which this is based is proportional to the number of vertices and the number of edges in the graph, assuming the jump function evaluations can be done in constant time. However, it is unrealistic to expect every jump function to be evaluated in constant time because, as we shall see, different strategies for constructing jump functions produce functions of varying execution costs.

Thus, we must argue about the number of times that a jump function is evaluated. A jump function J has $\text{support}(J)$ inputs, each of which can be lowered at most twice. Thus a jump function J can be evaluated no more than $O(|\text{support}(J)|)$ times. Let $\text{cost}(J)$ denote the cost of executing jump function J . For each jump function the total cost of execution will be $O(|\text{support}(J)| \cdot \text{cost}(J))$. Therefore the total cost of executing the interprocedural constant propagation algorithm is

$$(2.7) \quad O\left(\sum_s \sum_x |\text{support}(J_s^x)| \cdot \text{cost}(J_s^x)\right)$$

where s ranges over the call sites in the program and x ranges over the input parameters to the subroutine.

Construction of Jump Functions

Jump functions can vary widely in the precision and cost of their approximations. A jump function could involve full symbolic interpretation of the procedure it represents. On the other end of the spectrum, it could evaluate only the assignments that are on every path through the subroutine and not contained in any loop. In this case, variables assigned on optional control-flow paths or in loops would receive the value \perp (bottom) in the constant lattice.

The analyzer must construct jump functions for each procedure that has a call site in its body. Consider the example in Figure 13. To build a jump function for τ at call site γ , we need to know what action will be taken by subroutine INIT invoked at call site β . We address this problem by defining “return jump functions,” which summarize the constants propagated out of a subprogram when it is called with a particular set of inputs.

If x is an output of the procedure p , the *return jump function* R_p^x determines the value of x on return from an invocation of p in terms of the values of input parameters to p . The *support* of R_p^x is the same as the support of a forward jump function. In the simple case of the subroutine INIT shown in Figure 13 above, we have the following return jump function:

```

      PROGRAM MAIN
        INTEGER A
    α      CALL PROCESS(15,A)
        PRINT A
      END

      SUBROUTINE PROCESS(N,B)
        INTEGER N,B,I
    β      CALL INIT(I,N)
    γ      CALL SOLVE(B,I)
        RETURN
      END

      SUBROUTINE INIT(X,Y)
        INTEGER X,Y
        X = 2*Y
        RETURN
      END

      SUBROUTINE SOLVE(C,T)
        INTEGER C,T
        C = T*10
        RETURN
      END

```

FIGURE 13. A complex interprocedural constant folding example.

$$R_{\text{INIT}}^X = \{2*Y\}$$

and in the case of SOLVE, we have:

$$R_{\text{SOLVE}}^C = \{T*10\}$$

We can use return jump functions in the construction of forward jump functions to improve the results of constant propagation. For example, the jump function for call site γ is as follows

$$J_{\gamma}^T = \{\text{if } I \in \text{MOD}(\beta) \text{ then } R_{\text{INIT}}^X(N) \text{ else } \textit{undefined-const}\}$$

where *undefined-const* is used to signify a special value given to uninitialized variables. To round out the example, we present the remainder of the jump functions:

$$J_{\alpha}^N = \{15\} \text{ and } J_{\beta}^Y = \{N\}$$

One important jump function remains to be determined, namely the return jump function for subroutine `PROCESS`, which is the key to determining whether a constant can be substituted for the variable `A` in the print statement in the main program. This return jump function must invoke the return jump function for both `INIT` and `SOLVE` but the call to the return jump function for `INIT` will be automatic if we use the forward jump function to determine the value of the input formal parameter `T` of `SOLVE`. The resulting function is given below:

$$R_{\text{PROCESS}}^B = \begin{cases} \text{if } B \in \text{MOD}(\gamma) \\ \text{then } R_{\text{SOLVE}}^C(J_\gamma^T(N)) \\ \text{else } \textit{undefined-const} \end{cases}$$

Using these jump functions we can see that the value of the variable `A` on exit from the procedure `PROCESS` is given by

$$R_{\text{PROCESS}}^B = (2 * (15)) * 10 = 300$$

The study of scientific Fortran programs by Grove and Torczon [34] suggests that `MOD` should be computed prior to computing jump functions and prior to performing a global constant propagation to initialize the interprocedural propagation. Their results show that lack of `MOD` information stops propagation of constants within a procedure, and hence, within a jump function. Thus, `MOD` information improved the results of strictly intraprocedural constant propagation and the results of jump function evaluation during the interprocedural constant propagation. Since any variable that is not in `MOD` for a given call site has the identity function as its return jump function, using `MOD` information can simplify the forward jump functions. Of course, in a language where aliasing can occur, using `ALIAS` information in the analogous manner should help, too. The same study also established that:

1. increasing the complexity of the jump function implementations beyond some fairly simple approximations did not increase the set of constants found in scientific Fortran codes, and
2. return jump functions deliver new constants infrequently, but when they do, the payoff can be large.

In designing an interprocedural analyzer, the implementor must decide when to build jump functions. They can be constructed before analysis and optimization by a preliminary phase, as was done in the R^n programming environment [15, 29]. Alternatively, the analyzer can solve the simpler interprocedural problems first and use the results of `MOD` and `ALIAS` analysis during the construction of jump functions [34]. The former approach avoids some phase ordering difficulty and provides a clean separation between local analysis and interprocedural propagation. The jump functions built in this scheme explicitly test the values of `MOD` and `ALIAS` sets. The latter approach results in simpler jump functions, but requires the analyzer to examine each procedure between solving `MOD`

and ALIAS and performing constant propagation. The jump functions built in this scheme incorporate the MOD and ALIAS information implicitly, so they avoid explicit tests during propagation. (The examples shown earlier in this subsection assume the former approach.)

2.6. Kill Analysis

We can adapt some of the ideas used to solve the MOD and constant propagation problems to deal with kill analysis as well. For a block b , recall that $\text{KILL}(b)$ is the set of all variables that must be modified by execution of b . If the block is a call site s , then $\text{KILL}(s)$ is the set of all variables that must be modified on every path through the procedure called at s . For convenience in formulating the algorithm, we will compute $\text{NKILL}(s) = \neg\text{KILL}(s)$, the set of all variables that are not killed as a side effect of the call at s . A method similar to the one presented here can be used to compute $\text{USE}(s)$ for every call site s in the program.

Let us begin by considering how we would compute NKILL within a single procedure. Suppose we have, for each extended basic block b and each successor c of b , the set $\text{THRU}(b, c)$ of all variables that are not killed on some path through b to c . Then the following set of data-flow equations can be used to solve for $\text{NKILL}(b)$:

$$(2.8) \quad \text{NKILL}(b) = \bigcup_{c \in \text{succ}(b)} \text{THRU}(b, c) \cap \text{NKILL}(c)$$

If e is the exit node for the procedure then

$$(2.9) \quad \text{NKILL}(e) = \Omega$$

where Ω denotes the set of all variables. This set of equations can be solved using the simple iterative method of data-flow analysis.

To improve the precision of these NKILL sets, we can compute a set $\text{GNKILL}(p)$ for each procedure p in the program. Then, the THRU set for any block containing a call site s can be set to $\text{LOCAL} \cup \text{NKILL}(s)$, where $\text{NKILL}(s)$ is computed by projecting the GNKILL set for the called procedure back through the bindings that occur at s . This is more complicated than the single-procedure computation from Equation 2.8, because the solution for the procedures depend on each other. It is more complex than the MOD calculation because $\text{GNKILL}(p)$ depends on control-flow within p .

To model this intraprocedural control flow, we will construct a graph that we call the *reduced control-flow graph* $G_{\text{THRU}}(p)$ for procedure p . $G_{\text{THRU}}(p)$ is a graph in which each vertex is a call site, the entry node or the exit node for the procedure. Every edge (x, y) in G_{THRU} is annotated by the set $\text{THRU}(x, y)$ of variables that are not killed on some path from x to y not containing a call site. G_{THRU} is constructed by the algorithm in Figure 14. This algorithm is a

```

remove all back edges from the control-flow graph;
let  $b_0$  denote the procedure entry node;
mark  $b_0$  processed;
 $worklist := \emptyset$ ;
for each  $s \in successors(b_0)$  do  $worklist := worklist \cup \{(b_0, s)\}$  od;

while  $worklist \neq \emptyset$  do
  take an arbitrary element  $(b, s)$  from the worklist,
  such that all predecessors of  $s$  have already been processed
  or merged into  $b$ ;
  if  $s$  is a call site then
    for each  $t \in successors(s)$  do
       $worklist := worklist \cup \{(s, t)\}$  od;
    mark  $s$  as processed;
  else if  $s$  is the exit node then do nothing
  else /*  $s$  is normal node */
    merge  $s$  into  $b$ ;
    for each  $t \in successors(s)$  do
      if  $THRU[b, t]$  is undefined then
         $THRU[b, t] := THRU[b, s] \cap THRU[s, t]$ 
         $worklist := worklist \cup \{(b, t)\}$ 
      else
         $THRU[b, t] := THRU[b, t] \cup (THRU[b, s] \cap THRU[s, t])$ 
      fi
    od
  fi
od

```

FIGURE 14. Algorithm for constructing the reduced control-flow graph.

variant on topological sort; it is easy to see that it requires time linear in the size of the control flow graph.

Once we have the reduced control-flow graph, the algorithm shown in Figure 15 can be used to compute $GNKILL(p)$ for a procedure p given $GNKILL$ values for all procedures that can be called from within p .

A simple iterative data-flow analysis algorithm, employing the procedure in Figure 15, can be used to compute $GNKILL$, assuming that there are no call-by-reference formal parameters. Although this algorithm can take $O(N^2V)$ time in the worst case, if the call graph is reducible, the algorithm will require only $O((N + E)d)$ bit-vector steps, where d is the maximum number of back edges in any noncircular path in the call graph [46]. The total time required in this case is $O((N + E)dV)$. In practice, iterative methods converge quickly.

The algorithm described so far computes $GNKILL(p)$ rapidly for programs with no call-by-reference formal parameters. The algorithm can be adapted


```

for each  $b$  in  $G_{\text{THRU}}(p)$  in reverse topological order do
  if  $b$  is the exit node then  $\text{NKILL}[b] := \Omega$ 
  else if  $b$  is a call site then
     $\text{NKILL}[b] := \emptyset$ ;
    for each successor  $s$  of  $b$  in  $G_{\text{THRU}}(p)$  do
       $\text{NKILL}[b] := \text{NKILL}[b] \cup (\text{THRU}[b, s] \cap \text{NKILL}[s])$ 
    od;
     $\text{NKILL}[b] := \text{NKILL}[b] \cup \text{GNKILL}[q]$ 
    where  $q$  is the procedure called at  $b$ 
  else /*  $b$  is the entry node for  $G_{\text{THRU}}(p)$  */
     $\text{NKILL}[b] := \emptyset$ ;
    for each successor  $s$  of  $b$  in  $G_{\text{THRU}}(p)$  do
       $\text{NKILL}[b] := \text{NKILL}[b] \cup (\text{THRU}[b, s] \cap \text{NKILL}[s])$ 
    od;
     $\text{GNKILL}[p] := \text{NKILL}[b]$ 
  fi
od

```

FIGURE 15. Computing $\text{GNKILL}(p)$.

to produce the correct answers for programs with reference formal parameters if the appropriate actual-to-formal mappings are observed. However, it may take significantly longer to converge because of the possibility of the “shift-register effect” in which parameters are passed to other parameters in a recursive loop. This is the same problem we observed in dealing with the MOD problem in Section 2.3. Fortunately, we can use the same general technique—use of a formal parameter binding graph—to ameliorate this problem. We will construct a binding graph and mark it as shown in Figure 16.

This algorithm is simply a variation on the iterative algorithm, except that we update $\text{GNKILL}(p)$ for a procedure only if we discover that the status of one of its parameters might have changed to “killed” because a parameter that it is passed at some call site has changed its status. Since each parameter can be added to the worklist only once and since the algorithm visits each predecessor of a parameter taken from the worklist, the total number of GNKILL updates is limited to $O(E_B + N_B) = O(E + N)$, where E_B and N_B are the number of edges and vertices in the binding graph, respectively, and E and N are the number of edges and vertices in the call graph. The number of updates can be further reduced by carefully selecting the order of extraction of elements from the worklist.

Given GNKILL sets for every procedure in the program, we can derive a set $\text{NKILL}(s)$ for each call site s that invokes p by projecting $\text{GNKILL}(p)$ back through the bindings that occur at s . This is done with a formulation that is analogous to Equation 2.3:

initially, let the binding graph consist of a vertex
 for each formal parameter in the program;
 $worklist := \emptyset$;

for each procedure p in the program **do**
 let $GNKILL_0[p]$ be the result of applying the algorithm in
 Figure 15 with $GNKILL[q] = \Omega(q)$
 for each successor q of p ,
 where $\Omega(q)$ denotes the set of formal parameters of q ;
 $GNKILL[p] := GNKILL_0(p)$;
for each formal parameter f of p **do**
 if $f \in GNKILL_0[p]$ **then**
 $killed(f) := false$;
 for each formal parameter g to which f is passed
 at a call site within p **do**
 add an edge (f, g) to the binding graph
 od
 else
 $killed(f) := true$;
 $worklist := worklist \cup \{f\}$;
 fi
od
od;

while $worklist \neq \emptyset$ **do**
 select an arbitrary element $f \in worklist$;
 $worklist := worklist - \{f\}$;
for each g such that there is an edge (g, f) in the binding graph
 and $killed(g) = false$ **do**
 let q be the procedure of which g is a formal;
 $GNKILL[q] :=$ the result of applying the algorithm in
 Figure 15 with the current $GNKILL$ sets
 for its successors;
 if $g \notin GNKILL[q]$ **then**
 $killed(g) := true$;
 $worklist := worklist \cup \{g\}$;
 fi
od
od

FIGURE 16. Construct and mark the binding graph for GNKILL.

$$(2.10) \quad \text{NKILL}(s) = \{v \mid s \text{ invokes } p, v \xrightarrow{s} w \text{ and } w \in \text{GNKILL}(p)\}$$

Of course, the NKILL sets should also account for aliasing. What does it mean if, for some call site $s = (p, q)$, $v \in \text{NKILL}(s)$, $w \notin \text{NKILL}(s)$ and $w \in \text{ALIAS}(p, v)$? Because $v \in \text{NKILL}(s)$, we know that there is at least one path through q , and any procedures that it calls, along which v is not redefined. For w , the opposite is true; we know that every path through q , and any procedures that it calls, contains at least one redefinition of w . Because the ALIAS set contains MAY information, we do not know that v and w are aliases in every invocation of p . Thus, the sets are consistent. If we had MUST sets for ALIAS, the situation would be different. If $v \in \text{NKILL}(s)$, $w \notin \text{NKILL}(s)$ and $w \in \text{MUSTALIAS}(p, v)$, then we should remove v from $\text{NKILL}(s)$.

2.7. Advanced Analysis

In addition to the basic analysis problems discussed so far, successful program parallelization requires that some more sophisticated interprocedural problems be solved. Two that are especially important are symbolic analysis and array section analysis [44, 41].

Symbolic Analysis

Although constant propagation provides valuable information about variables on entry to a procedure, it is usually not possible to establish that a given variable is constant because, in most cases, it is not. However, it may not be necessary to prove a variable constant to improve global program analysis. It may be enough to establish bounds on its values or to prove that its value can be expressed as a function of the values of other variables at the same point in the program. Such symbolic relationships can be used to prove facts, such as the absence of a dependence. For example, consider the sample code below:

```
SUBROUTINE S(A,N,M)
  REAL A(N+M)
  INTEGER N, M
  DO I = 1, N
    A(I+M) = A(I) + B
  ENDDO
END
```

If we could prove that $N=M$ on entry to S , we could show that the loop within the subroutine carries no dependence.

The goal of *interprocedural symbolic analysis* is to prove facts about the values of variables that may hold on entry to a given subroutine or on return from a given call site. There are three types of symbolic analyses that are often carried out:

1. *symbolic expression analysis*, which seeks to determine a symbolic expression for the value of a variable in terms of the values of other variables at the same point in the program;
2. *predicate analysis*, which seeks to establish a relationship between values that a pair of variables may have at a given point in the program and
3. *range analysis*, which seeks to establish a range of values with known constant lower and upper bounds (and possibly strides) that a variable may take on at a given point in the program.

Often, the results of these analyses can be substituted for one another. For example, the fact that $M=N$ in the example above, could be established by either symbolic expression analysis or predicate analysis. The distinction between the two is that expression analysis can produce values that involve many other values, while predicate analysis typically relates only pairs of variables. On the other hand, the analysis that produces symbolic expressions can be more complicated than the analysis required for simple predicates.

Range analysis can be used effectively in program analysis to rule out certain possibilities. For example, consider the subroutine:

```

SUBROUTINE S(A,N,K)
  REAL A(0:N)
  INTEGER N, K
  DO I = 2, N
    A(I) = A(I) + A(K)
  ENDDO
END

```

If we can prove that $K \in [0:1]$ on entry to the subroutine, we can establish that the loop carries no dependence.

Within a single procedure, symbolic expression analysis is typically performed using some form of *value numbering* [4, 59], which uniquely numbers each expression so that expressions with equal value numbers at a given point will have equal values at run time. Since global program value numbering is likely to be complex, it is typically used only within a procedure. Interprocedurally, a restricted set of relationships, such as those represented by predicates involving only two variables, are propagated across procedure boundaries [43].

It is easy to see that range analysis and symbolic expression analysis can be handled by variations on the constant propagation algorithm from Section 2.5. To use this algorithm, we will need to define three things:

1. the process by which new symbolic information is introduced in a program;
2. *jump functions* which produce information at a call site from information at the entry to a procedure containing that call site; and
3. *return jump functions* which determine the relationships on output from a procedure given the relationships that hold on entry.

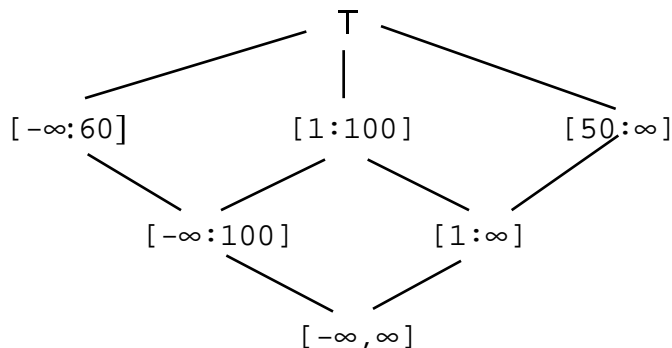


FIGURE 17. Simple value lattice for range analysis.

Let us consider how these three functions could be put together to compute range information. Range information is typically introduced at control-flow points in the program. For example, in a loop headed by

DO I = 1,N

it is safe to assume that $i \in [1:N]$. Similarly, conditional statements introduce partial ranges that can be composed to produce full ranges.

In all the symbolic analysis methods, jump functions compute values in lattices that are significantly more complex than the simple constant propagation lattice. For example, in the case of range analysis, we might use a lattice in which the meet operation picks the larger upper and smaller lower bound of the pair of ranges on the two joining control-flow branches. A fragment from such a lattice is depicted in Figure 17. If we bound the number of times that an upper bound can be increased before being taken as ∞ and we similarly bound the number of times a lower bound can be decreased, this lattice can be used anywhere a finite descending chain lattice is required, such as in an iterative algorithm.

Symbolic predicate analysis is more complicated than range analysis because it examines the relationship between a pair of variables. For example, it is useful to know whether two variables x and y are related by equality or are offset from one another by a constant. This might be represented as

$$x - y = c$$

where c is a compile-time constant. A more general linear relationship between variables can be characterized as:

$$c_1x + c_2y = c_0$$

where c_0 , c_1 and c_2 are all constants. Note that this relationship is transitive—if variables y and z are related by

$$d_1y + d_2z = d_0$$

then we can find constants e_0 , e_1 and e_2 such that

$$e_1x + e_2z = e_0$$

In particular,

$$c_1d_1x - c_2d_2z = c_0d_1 - c_2d_0$$

Thus, at any point in the program, we can find groups of variables that are linearly related to one another. The goal of symbolic predicate analysis is to propagate these sets throughout the program. This can also be done by employing a variant of interprocedural constant propagation.

Array Section Analysis

The analysis presented so far does not address one of the most important problems we need to solve if we are to automatically parallelize programs—namely, how to analyze dependences in loops that contain procedure calls. Consider the following code:

```

SUBROUTINE S
  DIMENSION A(100,100)
  ...
  DO I = 1,N
S1      CALL SOURCE(A,I) ! assigns to A
S2      CALL SINK(A,I)   ! uses A
  ENDDO
  RETURN
END

```

If we wish to parallelize the loop in this subroutine, we must determine whether any dependence is carried by the loop. Interprocedural information of the sort described in previous sections only tells us that array A is modified by SOURCE and used by SINK. Without better information, we must assume that there is an assignment in SOURCE to a location that is used by SINK on a later iteration of the loop—in other words, we must assume that the loop carries a dependence and cannot be parallelized.

The situation would be different if we were able to show that the accesses to A in both routines are confined to the i^{th} column, which is suggested by the use of I as a parameter to both routines. Then we know that different iterations of the loop deal with distinct portions of the array, so no carried dependence is possible. We would like to refine interprocedural analysis methods to be able to establish conditions like this, which means that our analysis needs to be able to recognize subarrays of the whole array.

Suppose that we are able to compute the set $M_A(i)$ of locations within the array that may be modified within SOURCE on the i^{th} iteration of the loop and the set $U_A(i)$ of locations that may be used in SINK on iteration i. These quantities

might be computed by array versions of MOD and USE respectively. Then the loop carries true dependence if and only if there exist indices I_1 and I_2 , $1 \leq I_1 < I_2 \leq N$, such that

$$M_A(I_1) \cap U_A(I_2) \neq \emptyset$$

In order to reason about subarrays, we need a method of representing them. The representation should be such that unions and intersections are reasonably easy to represent as well.

It is straightforward to extend the standard data-flow algorithms, which work on vectors of bits in which each bit can represent only two states (e.g., may be modified or must be preserved), to vectors of more general lattice elements. If we can find a lattice that represents subarrays accurately enough, we can use this lattice in our interprocedural data-flow analysis routines to determine side effects to subarrays.

Some important properties that a lattice representation should have are as follows:

1. the representation should be as accurate as possible;
2. the *meet* operation, which is invoked whenever two control-flow paths merge, must be efficient;
3. the dependence test, which usually involves intersection of region representations, should also be efficient;
4. it should be possible to handle recursion in the analysis framework, which implies that the lattice should have the *finite descending chain property*—that is, every descending chain in the lattice must reach a lattice minimum after a finite number of steps; and
5. it should be possible to deal with the parameter transformations that occur at call sites.

Let us consider one possible lattice for subarrays, depicted in Figure 18 below, that satisfies a number of these requirements. The elements of this lattice are referred to as *simple regular sections* because they can represent a very limited number of regular subarrays, namely points, rows, columns and the entire matrix.

Note that this lattice may extend to an infinite size because we can use arbitrary variables and constants in the subscripts. However, it has the finite descending chain property because no element in the lattice can be “lowered” more than three times.

Evaluating this lattice representation, we find that *meet*, which represents union in a case like the MOD calculation, has the following properties:

1. The depth of the lattice is $k+1$, where k is the number of subscript positions in the array represented.
2. The cost of a *meet* operation is $O(k)$, because each subscript position must be examined and compared for the two references to determine what the *meet* must be.

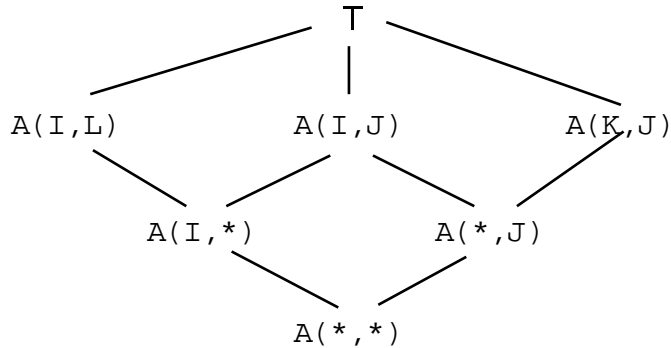


FIGURE 18. Simple regular section lattice.

3. Intersection, which is essential to the dependence test, is a limited form of unification, which can also be done in linear time in the number of subscripts.

To see this last claim, note that each subscript position in the lattice is either a symbolic expression or a constant or “*” which represents an entire row or column. Thus, if the two subscripts are symbolic expressions then the resulting subscript should be the same expression if they are equal and “*” otherwise. If one is an expression and the other is “*”, the resulting subscript should be “*”.

An important question remains: How accurate is this representation? It turns out that in practice this representation is too simple because it does not allow subarrays whose extents are bounded. Thus in an array $A(100,100)$ the best approximation to subarray $A(1:10,1:10)$ is $A(*,*)$. A much better representation is *bounded regular sections*, in which upper and lower bounds for each dimension are permitted. These can be thought of as any section represented by Fortran 90 triplet notation in which the stride is 1. Extensions that have been proposed in the literature include arbitrary stride triplet notation and triangular subarrays.

Interprocedural algorithms like the MOD solution can be directly adapted to deal with vectors of lattice elements, as long as the lattice has the finite descending chain property. The component of the MOD algorithm on the binding graph will converge because no formal parameter can be put on the worklist more than k times, where k is the maximum depth of the lattice. The reachability portion of the algorithm converges because, when a strongly-connected region is found, the element-vector for each element of the region is set to the minimum lattice element in the region.

Note that in the four side-effect problems described so far—MOD, REF, NKILL and USE—we always want to over-approximate the array section involved because we will optimize only when we *know* a variable cannot be in one of those sets for the call site in question. Recall the privatization example from the discussion of LIVE and USE in Section 2.1. We will be able to make a variable P private to a loop if we discover that, for the call site s_0 at the beginning of the loop:

$$P \in (\text{KILL}(S_0) \cap \neg \text{USE}(S_0)).$$

In other words, we optimize when

$$P \notin (\text{NKILL}(S_0) \cup \neg \text{USE}(S_0))$$

so we must overestimate rather than underestimate those sets to ensure that inaccuracy results only in missed opportunities, not in incorrect code.

2.8. Call Graph Construction

So far, we have assumed a precise call graph on which to solve interprocedural data-flow analysis problems. Constructing such a graph seems to be easy—simply examine each procedure in the program and, for each call site, construct an edge from the calling procedure to the called procedure. This simple approach works well, as long as there are no procedure parameters. Procedure parameters can add significant complications to the problem, as the following example shows:

```

SUBROUTINE SUB1(X,Y,P)
  INTEGER X, Y
S0    CALL P(X,Y)
      RETURN
      END

```

The problem here is determining what the called procedure might be at call site S_0 . If there is only one invocation sequence for SUB1 it may be simple to follow the call chain back to determine the procedure passed to parameter P . However, we cannot assume that a single chain exists, because procedure parameters were added to the language to ensure that many different procedures could be passed to P , even in the same program.

To solve this problem, we must be able to determine, for each procedure parameter P , the names of procedures that may be passed to P , directly or indirectly. However, we must be careful to avoid loss of precision in cases where more than one procedure parameter is passed. Consider the following example:

```

SUBROUTINE SUB2(X,P,Q)
  INTEGER X
S1    CALL P(X,Q)
      RETURN
      END

```

Suppose we have the following two calls:

```

CALL SUB2(X,P1,Q1)
CALL SUB2(X,P2,Q2)

```

where $P1$ and $P2$ simply invoke their procedure parameter on their integer parameter and return:

```

SUBROUTINE P1(X,Q)
  INTEGER X
S2    CALL Q(X)
      RETURN
END

```

Then at call site S_1 in subroutine SUB2, we can pass procedure Q1 to procedure P1 or we can pass procedure Q2 to procedure P2. We can *never* pass Q1 to P2 or Q2 to P1. In other words, P1 can only call Q1 and P2 can only call Q2 in this program. However a naive procedure-tracking scheme that simply maintains lists of procedures that could be passed to a given parameter might report edges from P1 to Q2 and from P2 to Q1 because the list for possible procedures passed to parameter Q in P1 includes both Q1 and Q2.

To overcome this problem, a precise call graph construction algorithm must keep track of which pairs of procedure parameters may be *simultaneously* passed to the procedure formal parameters in S_2 . This suggests a general algorithm for call graph construction.

Suppose we collect, for every procedure p that accepts procedure parameters, a set $\text{PROCPARMS}(p)$ of tuples of procedure names that may simultaneously be passed to p , where the order of the procedure names in the tuple corresponds to the order of the procedure parameters in the parameter list of p . Then the iterative algorithm in Figure 19 can be used to determine the correct set of procedure parameter tuples passed at each call site.

Theorem 5. For every procedure p in the program, the algorithm in Figure 19 produces $\text{PROCPARMS}(p)$ such that $\langle N_1, N_2, \dots, N_k \rangle \in \text{PROCPARMS}(p)$ if and only if there exists a call chain that passes parameter names N_1, N_2, \dots, N_k in that order, to the procedure parameter positions for p .

The proof of this theorem, which we do not include here, takes advantage of the fact that only a finite number of tuples of procedure calls can be in $\text{PROCPARMS}(p)$ for any given procedure p [61, 14]. The maximum number is given by the formula

$$(2.11) \quad \sum_p N_P^{\nu_p}$$

where ν_p is the number of procedure parameters to p .

We saw above that the number of tuples is given by the summation in Equation 2.11. Let $\nu_{\max} = \max_p(\nu_p)$, let N_C be the number of procedures with at least one procedure parameter, and let N_P be the number of procedure names that are passed to a procedure somewhere in the program. Then the running time can be approximated as

$$(2.12) \quad \sum_p N_P^{\nu_p} = O(N_C N_P^{\nu_{\max}}) \leq O(N N^{\nu_{\max}}) = O(N^{\nu_{\max}+1})$$

```

for each procedure  $p$  in the program do PROCPARMS[ $p$ ] :=  $\emptyset$  od;
 $W := \emptyset$ ;
for each call site  $s$  in the program do
  if the call site passes procedure names to all procedure
    parameters of the called procedure do
    let  $t = \langle N_1, N_2, \dots, N_k \rangle$  be the tuple of procedure names passed
      in order of the parameters to which they are passed;
     $W := W \cup \{ \langle t, p \rangle \}$ , where  $p$  is the procedure called
  fi
od;
while  $W \neq \emptyset$  do
  let  $\langle t = \langle N_1, N_2, \dots, N_k \rangle, p \rangle$  be an arbitrary element of  $W$ ;
   $W := W - \{ \langle t, p \rangle \}$ ;
  PROCPARMS[ $p$ ] := PROCPARMS[ $p$ ]  $\cup \{t\}$ ;
  let  $\langle P_1, P_2, \dots, P_k \rangle$  be the set of procedure parameters to which
    the elements of the tuple  $t = \langle N_1, N_2, \dots, N_k \rangle$  are mapped;
  for each call site  $s$  within  $p$  where  $P_i$  for some  $i, 1 \leq i \leq k$ ,
    is passed as a procedure parameter do
    let  $u = \langle M_1, M_2, \dots, M_k \rangle$  be the set of procedure names
      passed to the procedure  $q$  called at  $s$ , where each  $M_i$  is either
        the procedure name in the  $i$ th position or
         $N_j$  if  $P_j$  is passed in the  $i$ th position;
    if  $u \notin \text{PROCPARMS}[q]$  then  $W := W \cup \{ \langle u, q \rangle \}$  fi
  od
od

```

FIGURE 19. Algorithm for computing procedure parameter tuples.

where N is the number of procedures in the program. In the special case where there is no more than one procedure parameter to any procedure in the program, the running time is $O(N^2)$. In typical Fortran usage, the running time will not be a significant factor, because the use of procedure parameters is limited. However, for languages with more complex usage patterns, there exist approximate algorithms that run in linear or near-linear time in the size of the call graph [36, 38].

3. Interprocedural Optimization

3.1. Inline Substitution

The most familiar interprocedural optimization is inline substitution, by which the text of a subroutine is substituted at the point of call, with formal parameters replaced by actual parameters in the substituted text.

An example of inline substitution is presented below:

```

PROGRAM MAIN
  REAL A(100)
  CALL INPUT(A,N)
  DO I = 1,N
    CALL PROCESS(A,I)
  ENDDO
  CALL REPORT(A,N)
END
SUBROUTINE PROCESS(X,K)
  REAL X(*)
  X(K) = X(K) + K
  RETURN
END

```

If we inline subroutine PROCESS, substituting A and I for X and K respectively, we get the following code:

```

PROGRAM MAIN
  REAL A(100)
  CALL INPUT(A,N)
  DO I = 1,N
    A(I) = A(I) + I
  ENDDO
  CALL REPORT(A,N)
END

```

This code illustrates several well-known advantages of inline substitution:

1. procedure call overhead can be eliminated;
2. procedure body code can be tailored to the environment at the point of call, for example by using the index variable I from a register rather than memory; and
3. optimizations that would not be possible before substitution can be carried out—in this case, the loop can be vectorized.

The advantages of inlining are so compelling that many have suggested it as a generalized alternative to interprocedural analysis methods—if all the procedures in a program are inlined, ordinary single-processor analysis methods can be used to optimize the result.

However, overuse of inlining can cause a number of problems:

1. The massive substitution required can overwhelm the compilation system, since procedures after substitution might grow to unmanageable size, straining the capabilities of single-module compilers [45].
2. The object code generated from an inlined program may run slower because optimizing compilers do not handle codes resulting from systematic inlining well [25].

3. Any code that is changed inside an inlined subroutine will force the re-compilation of every procedure into which it has been substituted. In the limit, any change could require recompilation of the entire program.
4. Some subroutines are difficult to inline because of problems in substituting actual parameters for formals. This is illustrated by the example below:

```

PROGRAM MAIN
  REAL A(100, 100)
  ...
  CALL S(A(26,2),N)
  ...
END
SUBROUTINE S(X,M)
  REAL X(*)
  DO I = 1, M
    X(I) = X(I) + M
  ENDDO
  RETURN
END

```

The difficulty in this example arises because the two-dimensional actual parameter A is treated as a one-dimensional array in subroutine S. We could introduce an equivalence between array A in the main program and a single-dimensional array, to produce the following code:

```

PROGRAM MAIN
  REAL A(100, 100), a$(10000)
  EQUIVALENCE (A(1,1),a$(1))
  ...
  DO I = 1, N
    a$(I+125) = a$(I+125) + N
  ENDDO
  ...
END

```

The main problem with this approach is that it loses the information about the independence of different rows and columns, which can be critical to dependence analysis. Furthermore, this fix would not be available to us if the call to S were in another procedure which was passed A as a parameter.

Instead of systematic inlining, we recommend a selective, goal-directed inlining that uses global program analysis to determine when inlining would be profitable [7].

3.2. Procedure Cloning

Often the main benefit of inlining is the ability to take advantage of some specific optimizations that are possible at some but not every call site for that procedure. Consider the following example:

```

PROCEDURE UPDATE(A,N,IS)
  REAL A(N)
  INTEGER N, IS
  DO I = 1,N
    A(I*IS-IS+1) = A(I*IS-IS+1) + PI
  ENDDO
END

```

At first glance, this code looks vectorizable and it would be were it not for the possibility that the step size $IS = 0$, in which case the computation would become a sum of $N*PI$ into $A(1)$.

The obvious solution to this problem is to tailor the code to different specific versions based on the value of A . This could be done by a run-time test, but if we know the value of IS at compile time in every calling context, we can produce two versions of the program, one for the case $IS \neq 0$ and another for the case $IS = 0$. The compiler would be able to replace the call to `UPDATE` at the point of call with a call to one of the cloned procedures whenever it could determine the value at compile time.

Cloning is a particularly useful way to enhance the impact of constant propagation by treating parameters that are called with different constant values as constants in different clones of the original version of the procedure. This is the goal of cloning in the Convex Applications Compiler, the only commercial compiler we know that performs this optimization [56].

3.3. Hybrid Optimizations

There are cases where transformations involving more than one procedure can be used to gain some of the benefits of inlining without suffering the disadvantages. One such *hybrid optimization* is *loop embedding*, in which a loop is moved from one procedure to another [39].

Consider the original example for inlining in Section 3.1. If we interchanged the loop into subroutine `PROCESS`, we would get the following code inside that routine, which would vectorize well:

```

SUBROUTINE PROCESS(X,N)
  REAL X(*)
  DO K = 1,N
    X(K) = X(K) + K
  ENDDO
  RETURN
END

```

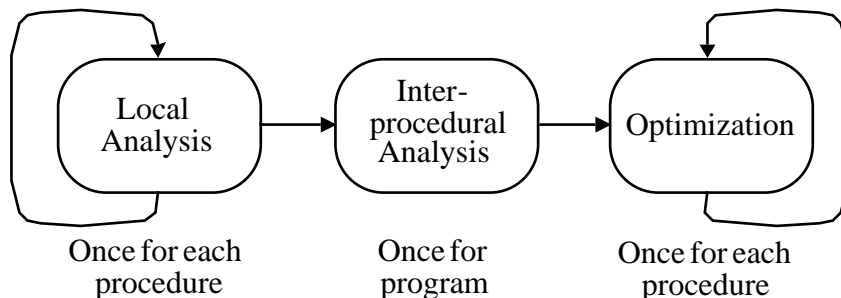


FIGURE 20. Interprocedural compilation process.

Note that the subroutine interface has changed to include the loop upper bound as a parameter instead of the loop index.

Interprocedural optimizations like this have been found useful in a number of cases, although there is as yet limited evidence as to their generality.

4. Managing Whole-Program Compilation

Using the results of interprocedural analysis and performing interprocedural optimizations complicates the compiler's job. Obviously, the compiler must perform additional tasks—at least, the analysis and optimization. More subtly, the use of these techniques introduces compilation dependences between the code produced for procedures of the compiled program. In a conventional compilation system, the object code for any single procedure is a function only of the source code for that procedure. In an interprocedural compilation system, the object code for a procedure may depend on the source code for the entire program. Thus, a change in source code could force recompilation of every procedure. To ensure correctness of the compiled code, the compiler must understand the compilation dependences that it introduces and ensure that they are respected. To satisfy users, the compiler must try to minimize the amount of recompilation performed in response to small changes in the source program.

We might expect that the interprocedural effects of changes made during the maintenance phase of a program's life would be somewhat limited. Thus, global program analysis methods that can examine the effects of interprocedural information flow might be useful in reducing the amount of recompilation.

We begin by subdividing the procedural compilation into two distinct phases, one that depends on interprocedural information and one that does not. The first of these phases, which we shall call *local analysis*, includes many of the usual compilation tasks—lexical analysis, parsing, and semantic analysis. At the same time, it examines the procedure for input to interprocedural analysis, determining the local sets used in each of the interprocedural analysis algorithms, such as the IMOD sets in MOD analysis. With this subdivision, the compilation process is structured as shown in Figure 20.

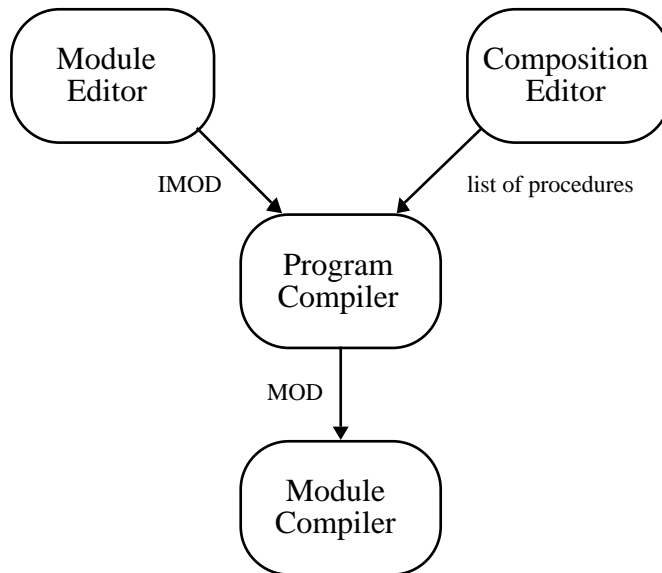
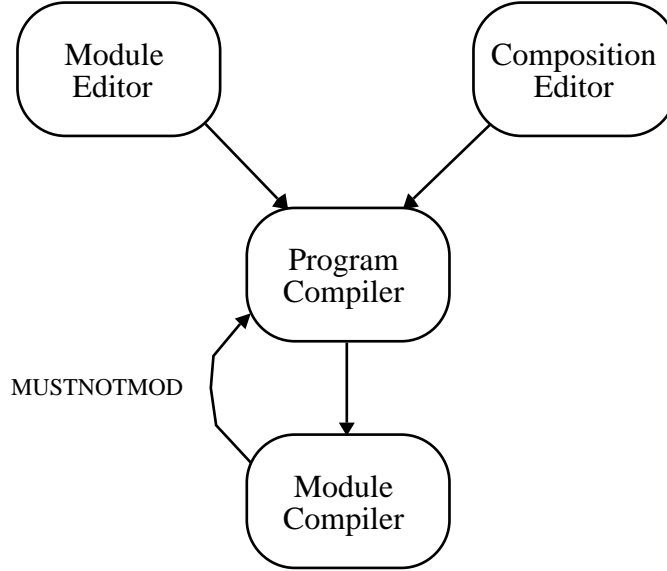


FIGURE 21. Interprocedural compilation system.

This organization permits interprocedural compilation, but does not solve the recompilation problem. However, if intermediate representations are saved, the local analysis phase will not need to be reinvoked for any unchanged procedures.

To allow us to conveniently compute and save intermediate information, we could organize our compilation system as shown in Figure 21. In this scheme, the job of lexical analysis is performed by a separate system component, which could reside in the module editor or in an importing tool. Information about each source procedure is stored in an intermediate representation that includes the parsed source. Programs are defined in this system by specification of a *program composition* via a *composition editor*, which can be viewed as an editor for lists of procedure names. The *program compiler* is the system component responsible for the compilation of the whole program. It reads all the local information for procedures in the program and carries out the various interprocedural analyses and optimizations required by the user. Once all interprocedural information is available, the *module compiler*, which is just a classic single-procedure optimizing compiler, optimizes each procedure and generates final code. Note that by separating the local analysis from the optimization, we have eliminated compilation-order dependences between the various procedures in the program.

This system, by itself, does not solve the recompilation problem. It must be coupled with a feed-back system from the individual module compilations that indicates which interprocedural analysis facts have been relied upon by the module compiler. To illustrate this process, consider the interprocedural recomputation of MOD sets, depicted in Figure 22.

**FIGURE 22.** Recompilation analysis.

We define the set $\text{MUSTNOTMOD}(p, s)$ to contain all the variables that must not be modified as a side effect of the procedure call at site s within p if the code generated for procedure p is to remain correct. Thus if the source code of procedure q called at s is changed so that a variable in $\text{MUSTNOTMOD}(p, s)$ is modified in q , then p must be recompiled. Think of $\text{MUSTNOTMOD}(p, s)$ as a record of the actions of the optimizing module compiler when it was last invoked on p . Whenever the optimizer used the fact that a variable was not modified at call site s to perform some optimization, it entered that variable into $\text{MUSTNOTMOD}(p, s)$.

In focusing on MUSTNOTMOD , we are relying on the observation that an optimizer will make program transformations based upon information only when it is sure that an undesirable event that would invalidate the optimization is not possible. Thus optimizations are based upon the *absence* of variables from $\text{MOD}(s)$ because that absence means that the variable cannot possibly be modified as a side effect of the call. Similarly, optimizations would be based on the absence of a variable from $\text{REF}(s)$, so $\text{MUSTNOTREF}(p, s)$ would be the corresponding recompilation set of interest. On the other hand, we will make optimizations like privatization based on the knowledge that a variable must be killed by a procedure call so $\text{MUSTKILL}(p, s)$ could be used to hold recompilation information.

In the case of MOD , let us consider what must be done by the program compiler to ensure correctness after a change to some procedure in the program.

1. First, the program compiler must recompute the MOD sets for every call site in the program.

2. For each call site $s \in p$ in the program, if

$$\text{MUSTNOTMOD}(p, s) \cap \text{MOD}(s) \neq \emptyset$$

then recompile procedure p .

The computation of MUSTNOTMOD sets is dependent on the optimizations used in the module compiler, which may be burdensome to collect accurately. However, simple approximation methods often work well in practice. Here are two such approximations for MUSTNOTMOD:

1. $\text{MUSTNOTMOD}(p, s) = \neg \text{MOD}(s)$ on the most recent compilation of p . Certainly, the module compiler cannot have depended on anything that was not true at the last compile.
2. $\text{MUSTNOTMOD}(p, s) = \neg \text{MOD}(s) \cap \text{REF}(p)$. This approximation goes one step farther than the previous one in that it does not force recompilation as a result of changes in status to variables that are merely passed through p .

The second of these approximations was used by the designers of the \mathbf{R}^n programming environment [17, 11].

5. Some Implementations

The authors have been personally involved with several interprocedural compilation systems:

1. The vectorization and parallelization system PFC [3], developed at Rice University, used interprocedural constant propagation and array side effect analysis to improve its dependence analysis. The results of dependence analysis could be displayed in a special browser called PTOOL. PFC was also used as a dependence server by later systems, including ParaScope described below.
2. The \mathbf{R}^n programming environment [17, 29], also developed at Rice, was the first one designed to support interprocedural analysis in a practical compilation system. Recompilation analysis, discussed in Section 4 was first tried in the \mathbf{R}^n system, which implemented constant propagation, ALIAS, MOD and REF.
3. ParaScope [22], a successor to \mathbf{R}^n , was designed to use interprocedural information in program parallelization. The FIAT interprocedural analysis framework [40] was originally developed for ParaScope, before being adapted for the SUIF compiler system (see below). Although it computed solutions to several interprocedural problems, ParaScope initially used PFC as a server for interprocedural array section analysis.

4. The D System [1], which was initially based on ParaScope, is intended to support programming in High Performance Fortran. It includes a framework that extends FIAT to perform interactive recompilation of interprocedural information when the user is editing in the context of a given program. In addition to classical analyses and optimizations, the D System propagates array distribution information interprocedurally to support compilation of Fortran D and High Performance Fortran [22].
5. The Stanford SUIF compiler was constructed with the aid of the interprocedural framework FIAT, a tool that facilitates rapid prototyping of interprocedural systems [40]. It performs a complete set of interprocedural analyses including constant propagation, kill analysis for arrays and live analysis. In a recent experiment, programs from the NAS, Spec-92 and Perfect benchmark suites were run through the compiler without any modification to the original source code. For the 27 programs in this test, 16 of them yielded more parallel loops as a result of interprocedural optimization. For 4 of these, significant speedups have been obtained only as a result of interprocedural optimization [37].
6. The Convex Application Compiler [56], which was modeled after the \mathbf{R}^n and ParaScope systems, is the first commercial system to perform interprocedural analysis on whole programs represented by multiple files. The Application Compiler does a complete job of interprocedural analysis and optimization for parallelization, computing the solutions to flow-sensitive interprocedural problems, cloning based on interprocedural constants, and employing array section analysis.

In addition to these, a number of commercial and research systems employ interprocedural information.

6. Historical Comments and References

Interprocedural analysis was introduced in a number of works published in the 1970s. Allen showed how interprocedural data-flow analysis could be carried out on programs without recursion [2]. In an unpublished abstract, Allen and Schwartz later extended these techniques to programs with recursion. Spillman described the interprocedural analysis that was available in the IBM PL/I compiler for procedures in a single file [66]. A principle goal of this implementation was analysis of pointer targets.

Barth was the first to publish a paper that discussed may and must problems [6]. Banning introduced the notion of flow-sensitive and flow-insensitive problems and presented a polynomial-time algorithm for solving such problems [5]. Myers presented a general algorithm for flow-insensitive problems, which he proved to be Co-NP Complete in the presence of aliasing [57].

Algorithms for flow-sensitive analysis have been described by Myers [57], Sharir and Pnueli [62], Harrison [42], Landi and Ryder [54], Choi, Burke and Carini [19], and Hall, Murphy and Amarasinghe [41]. The algorithm for flow-

insensitive interprocedural analysis presented in Section 2.3 is due to Cooper and Kennedy [26, 27] as is the alias analysis algorithm in Section 2.4 [28]. The flow-sensitive constant propagation algorithm presented in Section 2.5, is based on the work of Callahan, Cooper, Kennedy and Torczon [15]. The flow-sensitive kill analysis algorithm is modeled after an algorithm proposed by Callahan [13], although this particular formulation is new.

Interprocedural symbolic analysis has been discussed by a number of researchers, including Haghighat and Polychronopoulos [35], Irigoin, Jouvelot, and Triolet [50], and Hall, Murphy, and Amarasinghe [41]. The treatment in Section 2.7 follows Havlak [43].

Array section analysis, discussed in Section 2.7, has been the subject of work by a number of authors including Triolet, Irigoin and Feautrier [68], Burke and Cytron [10], Callahan and Kennedy [12, 16], Li and Yew [55], Havlak and Kennedy [44], Irigoin, Jouvelot, and Triolet [50], and Hind, Burke, Carini, and Midkiff [47]. Several researchers have developed algorithms for flow-sensitive array analysis, including Irigoin [49], Iitsuka [48], Tu and Padua [69], and Hall, Murphy and Amarasinghe [41].

Call graph analysis has been studied by Walter [70], Weihl [73], Spillman [66], Burke [9], and Shivers [63, 64, 65]. The call graph construction algorithm in Section 2.8 is based on an algorithm due to Ryder [61]. The proof that the same algorithm converges in the recursive case is due to Callahan, Carle, Hall and Kennedy [14]. An algorithm that is not precise but which generates its approximation in time linear in the size of the resulting call graph is given by Hall and Kennedy [36, 38].

Inline substitution has been widely studied [32, 58]. Cooper, Hall and Torczon noted its disadvantages [25]. Cloning was studied by Cooper, Hall and Kennedy [21, 23, 24]. Hybrid optimizations have been studied by Hall, Kennedy, and McKinley [39]. Techniques and algorithms closely related to cloning as described in this paper have been studied by Wegman [71] for intraprocedural analysis, Bulyonkov [8] and Ruf and Weise [60] for partial evaluation, Johnston [51] for dynamic compilation of APL, and Chambers and Ungar [18] for the language SELF.

The approach to program management described in Section 4 was pioneered in the \mathbf{R}^n programming environment [67, 29, 17]. Recompile analysis was developed by Burke, Cooper, Kennedy and Torczon [30, 11].

7. Summary

Although interprocedural analysis has limited applicability in optimizing code for uniprocessors, it is important for automatic parallelization systems. To be totally effective, a variety of analysis problems must be addressed. These include two types of problems:

1. forward propagation problems that determine the context in which a given procedure is called and

2. backward propagation problems which determine side effects of procedure calls.

In addition, interprocedural analysis problems can be classified as:

1. flow-insensitive problems, for which a precise solution does not require tracing control flow through the procedures down the call chain, and
2. flow-sensitive problems, which require control-flow tracing for a precise solution.

It has been shown that flow-insensitive forward and backward problems can be solved in $O((N + E)V)$ time, where N is the number of procedures in the program, E is the number of call sites and V is the number of global variables and parameters in the program.

In the most general case, flow-sensitive problems have been shown to be intractable. However, for typical programming languages, good approximations can be achieved in time that is polynomial in the size of the call graph. In practice, these algorithms perform in near linear time.

Interprocedural problems for single variables can be extended in a natural way to handle analysis of value ranges and symbolic values. In addition they can also be extended to analyze side effects to regular sections of array variables. In both cases, the running time will be expanded by a factor proportional to the maximum depth of the lattice of approximations used.

To be usable, an interprocedural analysis system needs to minimize the number of procedures that must be recompiled as a result of local program changes. Users will not tolerate a recompilation system that requires recompiling an entire program after a simple change to one procedure. The functions required to do this may be best embedded in a more general program management system.

Bibliography

- [1] V. Adve, A. Carle, E. Granston, S. Hiranandani, K. Kennedy, C. Koelbel, U. Kremer, J. M. Mellor-Crummey, C.-W. Tseng, and S. K. Warren. The d system: support for data-parallel programming. Technical Report CRPC-TR94-378, Center for Research on Parallel Computation, Rice University, Jan. 1994.
- [2] F. E. Allen. Interprocedural data flow analysis. In *Proceedings of the IFIP Congress 1974*, pages 398–402, Amsterdam, 1974. North Holland.
- [3] J. R. Allen and K. Kennedy. PFC: a program to convert FORTRAN to parallel form. In K. Hwang, editor, *Supercomputers: Design and Applications*, pages 186–203. IEEE Computer Society Press, Aug. 1984.
- [4] B. Alpern, M. N. Wegman, and F. K. Zadeck. Detecting equality of variables in programs. In *Conference Record of the Fifteenth ACM Symposium on the Principles of Programming Languages*, pages 1–11, San Diego, California, Jan. 1988.
- [5] J. P. Banning. An efficient way to find side effects of procedure calls and aliases of variables. In *Proceedings of the Sixth Annual Symposium on Principles of Programming Languages*, pages 29–41, San Antonio, Texas, Jan. 1979.
- [6] J. M. Barth. A practical interprocedural data flow analysis algorithm. *Communications of the ACM*, 21(9):724–736, Sept. 1978.

- [7] P. Briggs, K. Cooper, M. W. Hall, and L. Torczon. Goal-directed interprocedural optimization. Technical Report TR90-148, Dept. of Computer Science, Rice University, Nov. 1990.
- [8] M. A. Bulyonkov. Polyvariant mixed computation for analyzer programs. *Acta Informatica*, 21:473–484, 1984.
- [9] M. Burke. An interval-based approach to exhaustive and incremental interprocedural data-flow analysis. *ACM Trans. Prog. Lang. Syst.*, 12(3):341–395, July 1990.
- [10] M. Burke and R. Cytron. Interprocedural dependence analysis and parallelization. *SIGPLAN Notices*, 21(7):162–175, July 1986. *Proceedings of the ACM SIGPLAN 86 Symposium on Compiler Construction*.
- [11] M. Burke and L. Torczon. Interprocedural optimization: Eliminating unnecessary recompilation. *ACM Trans. Prog. Lang. Syst.*, 15(3):367–399, July 1993.
- [12] D. Callahan. *A global approach to the detection of parallelism*. PhD thesis, Department of Computer Science, Rice University, Houston, TX, Dec. 1987.
- [13] D. Callahan. The program summary graph and flow-sensitive interprocedural data flow analysis. *SIGPLAN Notices*, 23(7):47–56, July 1988. *Proceedings of the ACM SIGPLAN '88 Conference on Programming Language Design and Implementation*.
- [14] D. Callahan, A. Carle, M. W. Hall, and K. Kennedy. Constructing the procedure call multigraph. *IEEE Transactions on Software Engineering*, 16(4):483–487, Apr. 1990.
- [15] D. Callahan, K. D. Cooper, K. Kennedy, and L. Torczon. Interprocedural constant propagation. *SIGPLAN Notices*, 21(7):152–161, July 1986. *Proceedings of the ACM SIGPLAN '86 Symposium on Compiler Construction*.
- [16] D. Callahan and K. Kennedy. Analysis of interprocedural side effects in a parallel programming environment. *Journal of Parallel and Distributed Computing*, 5:517–550, 1988.
- [17] A. Carle, K. D. Cooper, R. T. Hood, K. Kennedy, L. Torczon, and S. K. Warren. A practical environment for Fortran programming. *IEEE Computer*, 20(11):75–89, Nov. 1987.
- [18] C. Chambers and D. Ungar. Customization: Optimizing compiler technology for SELF, a dynamically-typed object-oriented programming language. *SIGPLAN Notices*, 24(7):146–160, July 1989. *Proceedings of the ACM SIGPLAN '90 Conference on Programming Language Design and Implementation*.
- [19] J.-D. Choi, M. Burke, and P. Carini. Efficient flow-sensitive interprocedural computation of pointer-induced aliases and side effects. In *Conference Record of the Twentieth Annual ACM Symposium on Principles of Programming Languages*, pages 232–245, Charleston, South Carolina, Jan. 1993.
- [20] J. Cocke and P. Markstein. Measurement of program improvement algorithms. In *Proceedings of Information Processing 80*, pages 221–228. North Holland Publishing Company, 1980.
- [21] K. D. Cooper. *Interprocedural Data Flow Analysis in a Programming Environment*. PhD thesis, Department of Mathematical Sciences, Rice University, Houston, Texas, Apr. 1983.
- [22] K. D. Cooper, M. Hall, R. T. Hood, K. Kennedy, K. McKinley, J. M. Mellor-Crummey, L. Torczon, and S. K. Warren. The ParaScope parallel programming environment. *Proceedings of the IEEE*, 81(2):244–263, Feb. 1993.
- [23] K. D. Cooper, M. W. Hall, and K. Kennedy. Procedure cloning. In *Proceedings of the IEEE International Conference on Computer Languages*, pages 96–105. IEEE, Apr. 1992.
- [24] K. D. Cooper, M. W. Hall, and K. Kennedy. A methodology for procedure cloning. *Computer Languages*, 19(2):105–118, Apr. 1993.
- [25] K. D. Cooper, M. W. Hall, and L. Torczon. An experiment with inline substitution. *Software-Practice and Experience*, 21(6):581–601, June 1991.
- [26] K. D. Cooper and K. Kennedy. Efficient computation of flow-insensitive interprocedural summary information—a correction. *SIGPLAN Notices*, 23(4):35–42, Apr. 1988.
- [27] K. D. Cooper and K. Kennedy. Interprocedural side-effect analysis in linear time. *SIGPLAN Notices*, 23(7):57–66, July 1988. *Proceedings of the ACM SIGPLAN '88 Conference on Programming Language Design and Implementation*.

- [28] K. D. Cooper and K. Kennedy. Fast interprocedural alias analysis. In *Conference Record of the Sixteenth Annual ACM Symposium on Principles of Programming Languages*, pages 49–59, Austin, Texas, Jan. 1989.
- [29] K. D. Cooper, K. Kennedy, and L. Torczon. The impact of interprocedural analysis and optimization in the r^n programming environment. *ACM Trans. Prog. Lang. Syst.*, 8(4):491–523, Oct. 1986.
- [30] K. D. Cooper, K. Kennedy, and L. Torczon. Interprocedural optimization: Eliminating unnecessary recompilation. *SIGPLAN Notices*, 21(7):58–67, July 1986. *Proceedings of the ACM SIGPLAN '86 Symposium on Compiler Construction*.
- [31] 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 Trans. Prog. Lang. Syst.*, 13(4):451–490, Oct. 1991.
- [32] J. W. Davidson and A. M. Holler. A study of a c function inliner. *Software-Practice and Experience*, 18(8):775–790, Aug. 1988.
- [33] S. L. Graham and M. N. Wegman. A fast and usually linear algorithm for global data flow analysis. *Journal of the ACM*, 23(1):172–202, 1976.
- [34] D. Grove and L. Torczon. Interprocedural constant propagation: A study of jump function implementations. *SIGPLAN Notices*, 28(6):90–99, June 1993. *Proceedings of the ACM SIGPLAN '93 Conference on Programming Language Design and Implementation*.
- [35] M. Haghighat and C. Polychronopoulos. Symbolic analysis: A basis for parallelization, optimization, and scheduling of programs. In *Proceedings of the Sixth Workshop on Languages and Compilers for Parallel Computing*, Portland, OR, Aug. 1993.
- [36] M. W. Hall. *Managing Interprocedural Optimization*. PhD thesis, Department of Computer Science, Rice University, Houston, Texas, Apr. 1991.
- [37] M. W. Hall, S. Amarasinghe, B. Murphy, and M. S. Lam. Interprocedural analysis for parallelization: preliminary results. Technical Report CSL-TR-95-665, Stanford Computer Systems Laboratory, Mar. 1995.
- [38] M. W. Hall and K. Kennedy. Efficient call graph analysis. *ACM Letters on Programming Languages and Systems*, 1(3):227–242, Sept. 1992.
- [39] M. W. Hall, K. Kennedy, and K. S. McKinley. Interprocedural transformations for parallel code generation. In *Proceedings of Supercomputing '91*, Albuquerque, NM, Nov. 1991. IEEE Computer Society Press.
- [40] M. W. Hall, J. M. Mellor-Crummey, A. Carle, and R. G. Rodriguez. Fiat: a framework for interprocedural analysis and transformation. In *Proceedings of the Sixth Workshop on Languages and Compilers for Parallel Computing*, Portland, OR, Aug. 1993.
- [41] M. W. Hall, B. Murphy, and S. Amarasinghe. Interprocedural analysis for parallelization: A case study. In *Proceedings of the Seventh SIAM Conference on Parallel Processing for Scientific Computing*, San Francisco, CA, Feb. 1995.
- [42] W. L. Harrison. The interprocedural analysis and automatic parallelization of scheme programs. *Lisp and Symbolic Computation*, 2(3/4):179–396, Oct. 1989.
- [43] P. Havlak. *Interprocedural Symbolic Analysis*. PhD thesis, Department of Computer Science, Rice University, Houston, Texas, May 1994.
- [44] P. Havlak and K. Kennedy. An implementation of interprocedural bounded regular section analysis. *IEEE Transactions on Parallel and Distributed Systems*, 2(3):350–360, July 1991.
- [45] M. S. Hecht. *Flow Analysis of Computer Programs*. Elsevier North-Holland, 1977.
- [46] M. S. Hecht and J. D. Ullman. A simple algorithm for global data flow analysis of programs. *SIAM Journal of Computing*, 4:519–532, 1975.
- [47] M. Hind, M. Burke, P. Carini, and S. Midkiff. An empirical study of precise interprocedural array analysis. *Scientific Programming*, 3(3):255–271, 1994.
- [48] T. Iitsuka. Flow-sensitive interprocedural analysis method for parallelization. In *IFIP TC10/WG10.3 Working Conference on Architectures and Compilation Techniques for Fine and Medium Grain Parallelism*, Jan. 1993.
- [49] F. Irigoin. Interprocedural analyses for programming environments. In *NSF-CNRS Work-*

- shop on Environments and Tools for Parallel Scientific Programming*, Sept. 1992.
- [50] F. Irigoin, P. Jouvelot, and R. Triolet. Semantical interprocedural parallelization: An overview of the pips project. In *Proceedings of the 1991 ACM International Conference on Supercomputing*, Cologne, Germany, June 1991.
 - [51] R. Johnston. The dynamic incremental compiler of apl \ 3000. In *Proceedings of the APL '79 Conference*, pages 82–87. ACM, June 1979.
 - [52] J. Kam and J. D. Ullman. Global data flow analysis and iterative algorithms. *Journal of the ACM*, 23(1):159–171, Jan. 1976.
 - [53] J. Kam and J. D. Ullman. Monotone data flow analysis frameworks. *Acta Informatica*, 7(3):305–318, 1977.
 - [54] W. Landi and B. G. Ryder. A safe approximate algorithm for interprocedural pointer aliasing. *SIGPLAN Notices*, 27(7):235–248, July 1992. *Proceedings of the ACM SIGPLAN '92 Conference on Programming Language Design and Implementation*.
 - [55] Z. Li and P.-C. Yew. Efficient interprocedural analysis for program parallelization and restructuring. *SIGPLAN Notices*, pages 85–99, July 1988. *Proceedings of the ACM SIGPLAN Symposium on Parallel Programming: Experience with Applications, Languages, and Systems (PPEALS)*.
 - [56] R. Metzger and S. Stroud. Interprocedural constant propagation: An empirical study. *ACM Letters on Programming Languages and Systems*, 2(1–4):213–232, March–December 1993.
 - [57] E. W. Myers. A precise interprocedural data flow algorithm. In *Conference Record of the Eighth Annual ACM Symposium on Principles of Programming Languages*, pages 219–230, Williamsburg, Virginia, Jan. 1981.
 - [58] S. Richardson and M. Ganapathi. Interprocedural analysis versus procedure integration. *Information Processing Letters*, 32(1), Aug. 1989.
 - [59] B. K. Rosen, M. N. Wegman, and F. K. Zadeck. Global value numbers and redundant computations. In *Conference Record of the Fifteenth Annual ACM Symposium on Principles of Programming Languages*, pages 12–27, San Diego, California, Jan. 1988.
 - [60] E. Ruf and D. Weise. Using types to avoid redundant specialization. *SIGPLAN Notices*, 26(9):321–333, Sept. 1991. *Proceedings of the PEPM '91 Symposium on Partial Evaluation and Semantics-Based Program Manipulation*.
 - [61] B. G. Ryder. Constructing the call graph of a program. *IEEE Transactions on Software Engineering*, SE-5(3):216–225, 1979.
 - [62] 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*. Prentice Hall, Englewood Cliffs, NJ, 1981.
 - [63] O. Shivers. Control-flow analysis in Scheme. *SIGPLAN Notices*, 23(7):164–174, July 1988. *Proceedings of the ACM SIGPLAN '88 Conference on Programming Language Design and Implementation*.
 - [64] O. Shivers. *Control-Flow Analysis of Higher-Order Languages*. PhD thesis, School of Computer Science, Carnegie Mellon University, Pittsburg, PA, May 1991.
 - [65] O. Shivers. The semantics of scheme control flow analysis. *SIGPLAN Notices*, 26(9):190–198, Sept. 1991. *Proceedings of the Symposium on Partial Evaluation and Semantics-Based Program Manipulation*.
 - [66] T. C. Spillman. Exposing side-effects in a PL/I optimizing compiler. In *Proceedings of the IFIP Congress 1971*, pages 376–381. North Holland, 1971.
 - [67] L. Torczon. *Compilation Dependences in an Ambitious Optimizing Compiler*. PhD thesis, Department of Computer Science, Rice University, Houston, Texas, Apr. 1985.
 - [68] R. Triolet, F. Irigoin, and P. Feautrier. Direct parallelization of call statements. *SIGPLAN Notices*, 21(7):176–185, July 1986. *Proceedings of the ACM SIGPLAN '86 Symposium on Compiler Construction*.
 - [69] P. Tu and D. Padua. Automatic array privatization. In *Proceedings of the Sixth Workshop on Languages and Compilers for Parallel Computing*, Portland, OR, Aug. 1993.
 - [70] K. Walter. Recursion analysis for compiler optimization. *Communications of the ACM*,

- 19(9):514–516, 1976.
- [71] M. N. Wegman. *General and Efficient Methods for Global Code Improvement*. PhD thesis, University of California, Berkeley, CA, Dec. 1981.
 - [72] M. N. Wegman and F. K. Zadeck. Constant propagation with conditional branches. *ACM Trans. Prog. Lang. Syst.*, 13(2):181–210, Apr. 1991.
 - [73] W. E. Weihl. Interprocedural data flow analysis in the presence of pointers, procedure variables and label variables. In *Conference Record of the Seventh Annual ACM Symposium on Principles of Programming Languages*, pages 83–94, Las Vegas, Nevada, Jan. 1980.

Received February 1995.