

Second-Order Abstract Interpretation via Kleene Algebra

Łucja Kot

lucja@cs.cornell.edu

Dexter Kozen

kozen@cs.cornell.edu

Department of Computer Science

Cornell University

Ithaca, New York 14853-7501, USA

Abstract

Most standard approaches to the static analysis of programs, such as the popular worklist method, are first-order methods that inductively annotate program points with abstract values. In this paper we introduce a second-order approach based on Kleene algebra. In this approach, the primary objects of interest are not the abstract data values, but the transfer functions that manipulate them. These elements form a left-handed Kleene algebra. The dataflow labeling is not achieved by inductively labeling the program with abstract values, but rather by computing the star (Kleene closure) of a matrix of transfer functions. In this paper we introduce the method and prove soundness and completeness with respect to the standard worklist algorithm.

1 Introduction

Dataflow analysis and abstract interpretation are concerned with the static derivation of information about the execution state at various points in a program. There is typically a semilattice L of *types* or *abstract values*, each describing a larger set of possible runtime values. The natural partial order associated with L is the subtype relation. The objective of the analysis is to associate an element of L with each point of the program that represents what is known about the program state whenever control passes through that point. This may consist of type information, bounds on values of variables or registers, operand stack depth, the shape of data structures, whether pointers are null, etc. It is usually only an approximation; the lower in the semilattice, the better the approximation.

Each instruction has one or more associated *transfer functions* $f : L \rightarrow L$ that describe how the state is transformed by the instruction. The domain of f is determined by the type of the instruction. Membership in the domain of f may be considered a precondition for the safe execution

of the instruction; attempting to apply f to an element of L not in its domain signals a type error. For example, an empty stack should not be popped. Transfer functions may be composed, provided there is no type mismatch. Transfer functions may be polymorphic, as for example in the case of a swap operation that interchanges the top two elements of a stack.

Each instruction has zero or more *successors*. Data manipulation instructions such as loads, stores, and arithmetic operations have the fallthrough instruction as successor. Conditional jumps have the fallthrough as well as the jump target, and unconditional jumps have only the jump target. Any instruction that can raise an exception has the entry point of an exception handler as a successor. Successors are determined solely by statically available information. Thus the program can be modeled by a directed control flow graph G whose nodes are the instructions and whose edges go from each instruction to the successors of that instruction. The edges are labeled with the transfer functions associated with that instruction. In most cases the transfer function is the same for all edges exiting a particular node, but in some cases it is different. For example, in Java bytecode, if an instruction throws an exception, then the operand stack is cleared and the exception object pushed onto the stack before invoking the handler. The successor state corresponding to the exception handler thus reflects a different stack configuration than that corresponding to the fallthrough instruction.

The *worklist algorithm* for dataflow analysis is a standard method for computing a least fixpoint labeling of the nodes of G with elements of L [6]. It works as follows. First, the entry point of the method is labeled with the element of L describing the initial state of the computation. For example, in Java bytecode, the initial label consists of an empty operand stack, the types of the arguments to the method (including the object itself if it is an instance method) in the first few local variables, and a special undefined marker for the remaining local variables. This node

is marked as changed and placed on a worklist. Then, as long as the worklist is nonempty, the procedure repeatedly removes the next element s of the worklist, and for each exiting edge (s, t) , applies the transfer function f associated with that edge to the label $x \in L$ of s to get $f(x)$, then updates the label of the successor t with $f(x)$. If t is unlabeled, then it is labeled with $f(x)$. If t is already labeled, then it is relabeled with the join of $f(x)$ and its current label in the semilattice, if the join exists. This indicates the best information that is known about the program state at t from the various control flow paths into t that have been analyzed so far. If the join of $f(x)$ and the current label of t does not exist in the semilattice, then it is a type error. If t is successfully labeled and the new label is different from the old, then t is marked as changed and placed back on the worklist. When the worklist becomes empty, the resulting labeling is the least fixpoint of a monotone mapping on labelings defined in terms of the transfer functions.

One disadvantage of the worklist approach is that long paths in the graph may be analyzed several times. For example, if a node s is labeled with $x \in L$, then later revisited and relabeled with $y > x$, then any long paths out of s may be traversed again. The running time could be as bad as dn , where n is the size of the program and d is the depth of the semilattice, although this worst-case bound is probably rarely attained in practice. Thus the worklist algorithm remains a popular method for many practical program analysis tasks.

In this paper we describe an alternative approach that can be used to avoid the recalculation of dataflow information along long paths using a symbolic method based on Kleene algebra. The elements of the algebra are transfer functions. The novelty of this approach is that it is the transfer functions, not the data values, that are the objects of primary algebraic interest. The transfer functions are elements of a certain algebraic structure called a *left-handed Kleene algebra* with the operations of composition, join, and iteration. The control flow graph of a method with n instructions gives rise to an $n \times n$ matrix of transfer functions, and computing the star or Kleene closure of this matrix amounts to computing the dataflow information at all points of the program simultaneously.

In a companion paper [8], we describe a concrete application of these ideas in the context of Java bytecode. That application combines the second-order approach introduced in this paper with the standard worklist algorithm to obtain a hybrid algorithm with running time $O(nm + m^3)$, where m is the size of a cutset (a set of nodes breaking all cycles in G). The algorithm avoids recalculation of dataflow information along long paths by computing the star (closure) of an $m \times m$ matrix of transfer functions. This may give an improvement when m is small compared to n .

In this paper, we lay the foundations of this approach

and prove correctness with respect to the standard worklist algorithm [6]. This paper is organized as follows. In Section 2, we review the pertinent definitions of Kleene algebra and left-handed Kleene algebra, including the formation of matrices, and describe how transfer functions can be modeled as strict monotone functions on a semilattice of abstract values or types. In Section 3, we present an alternative approach to static analysis based on computing the star or Kleene closure of a matrix of transfer functions. Finally, in Section 4 we prove that our second-order method produces the same final dataflow labeling as the standard worklist algorithm on all type-correct programs.

1.1 KAT and the Static Analysis of Programs

In [10], we showed how Kleene algebra with tests (KAT), a variant of KA that includes Boolean tests, can be used to statically verify compliance with safety policies specified by *security automata*, a general mechanism for the specification and enforcement of a large class of safety policies [12]. We proved the soundness and completeness of the method over relational interpretations, and illustrated the method on an example of [3] involving the verification of a device driver.

The results of [10] are very different from those of the present paper. In that paper, the objective was to show how the deductive system could be used as a mechanism to propagate state information throughout the program. That was also a first-order approach. The Boolean algebra and the deductive system were essential components of that program. Here, we are not restricting ourselves to Boolean information, and we are not using the deductive system.

2 Background

2.1 Upper Semilattices

An *upper semilattice* is a partially ordered set L in which every finite set has a least upper bound, which must be unique. The least upper bound of two elements x and y is denoted $x + y$. The least upper bound of the null set is denoted \perp . The operation $+$ is associative, commutative, idempotent ($x + x = x$), and $x \leq y$ iff $x + y = y$. The element \perp is the least element of the semilattice and is an identity for $+$.

We also assume that there are no infinite ascending chains in L ; this is known as the *ascending chain condition* (ACC). It follows from this assumption that there exists a maximum element \top .

Intuitively, lower elements in the semilattice represent more specific information, and the join operation represents disjunction of information. For example, in the Java

class hierarchy, the join of `String` and `StringBuffer` is `Object`, their least common ancestor in the hierarchy.

The element \top represents a type error. In practice, any attempt by a dataflow analysis computation to form a join $x + y$ that does not make sense indicates a fatal type error, and the analysis will be aborted. We represent this situation mathematically by $x + y = \top$.

The element \perp represents “unlabeled”. For example, the initial labeling in the worklist algorithm is a map $w_0 : V \rightarrow L$, where V is the set of vertices of the control flow graph, such that $w_0(s_0)$ is the initial dataflow information available at the start node s_0 , and $w_0(u) = \perp$ for all other nodes $u \in V$.

The ascending chain condition (ACC) is a standard assumption that ensure that dataflow computations always converge.

2.2 Kleene Algebra

Kleene algebra (KA) was introduced by S. C. Kleene [7] (see also [4]). We define a *Kleene algebra* (KA) to be a structure $(K, +, \cdot, *, 0, 1)$, where $(K, +, \cdot, 0, 1)$ is an idempotent semiring, a^*b is the least x such that $b+ax \leq x$, and ba^* the least x such that $b + xa \leq x$. Here “least” refers to the natural partial order $a \leq b \leftrightarrow a + b = b$. The operation $+$ gives the supremum with respect to \leq . This particular axiomatization is from [9]. We normally omit the \cdot , writing ab for $a \cdot b$. The precedence of the operators is $* > \cdot > +$. Thus $a + bc^*$ should be parsed $a + (bc^*)$.

KA has a rich algebraic theory with many natural and useful models: language-theoretic, relational, trace-based, matrix. Standard models include the family of regular sets of strings over a finite alphabet, the family of binary relations on a set, and the family of $n \times n$ matrices over another Kleene algebra. We refer the reader to [9] for further definitions and basic results.

For this paper, we consider a weaker axiomatization. We will assume that the algebra is left-distributive, but not necessarily right-distributive. However, we will assume that it is right-*predistributive*. That is, we assume that $ab + ac = a(b + c)$, but only $ac + bc \leq (a + b)c$. Moreover, we will not require that a^*b be the least x such that $b + ax \leq x$, but only that ba^* be the least x such that $b + xa \leq x$. Such algebras are called a *left-handed Kleene algebras*.

The following summarizes the axioms of left-handed KA:

$$\begin{array}{ll} a + (b + c) = (a + b) + c & a(bc) = (ab)c \\ a + b = b + a & 1a = a1 = a \\ a + 0 = a + a = a & 0a = a0 = 0 \\ ab + ac = a(b + c) & ac + bc \leq (a + b)c \end{array}$$

and for the $*$ operator,

$$1 + a^*a \leq a^* \quad (1)$$

$$b + xa \leq x \Rightarrow ba^* \leq x, \quad (2)$$

where \leq refers to the natural partial order on K :

$$a \leq b \Leftrightarrow a + b = b.$$

Instead of (2), we might take the equivalent axiom

$$xa \leq x \Rightarrow xa^* \leq x \quad (3)$$

One can show that all the operations are monotone with respect to \leq . The proof of monotonicity of multiplication does not need distributivity, but only predistributivity. One can also show that the inequality (1) can be strengthened to an equality.

2.3 Matrices

As mentioned, the $n \times n$ matrices over a Kleene algebra again form a Kleene algebra under the appropriate definitions of the operators. One can establish that matrices over a left-handed algebra are left-handed, and those over a right-handed algebra are right-handed. The proofs are symmetric. Unfortunately, the proof given in [9] uses distributivity on both sides to show only right-handedness, so technically it does not suffice to establish the result we need. We therefore supply a proof here for completeness. We show that 2×2 matrices over a left-handed KA are left-handed; the result for general n follows by induction as in [9]. We also show only (1) and (3); the other laws are all straightforward.

Define the matrix E^* from E as follows:

$$E = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad E^* = \begin{bmatrix} f^* & f^*bd^* \\ g^*ca^* & g^* \end{bmatrix},$$

where $f = a + bd^*c$ and $g = d + ca^*b$. The inequality (1) for E is $I + E^*E \leq E^*$, which reduces to the following inequalities over K :

$$\begin{array}{l} 1 + f^*a + f^*bd^*c \leq f^* \\ f^*b + f^*bd^*d \leq f^*bd^* \\ g^*ca^*a + g^*c \leq g^*ca^* \\ 1 + g^*ca^*b + g^*d \leq g^*. \end{array}$$

For the first two,

$$\begin{aligned} 1 + f^*a + f^*bd^*c &= 1 + f^*(a + bd^*c) \\ &= 1 + f^*f \leq f^*, \\ f^*b + f^*bd^*d &= f^*b(1 + d^*d) \leq f^*bd^*, \end{aligned}$$

and the other two are symmetric. To show (3) for E , we must show that $XE \leq X$ implies $XE^* \leq X$. We can

show this independently for each row of X . This reduces to the task of showing

$$xf^* + yg^*ca^* \leq x \quad (4)$$

$$xf^*bd^* + yg^* \leq y \quad (5)$$

under the assumptions

$$xa + yc \leq x \quad xb + yd \leq y.$$

By symmetry, we need only show (4). By the property (3) for K , we have $xa^* \leq x$ and $yd^* \leq y$, therefore

$$\begin{aligned} xf &= x(a + bd^*c) \leq xa + xbd^*c \\ &\leq xa + yd^*c \leq xa + yc \leq x. \end{aligned}$$

It follows from (3) that $xf^* \leq x$. By a symmetric argument, $yg^* \leq y$. Thus

$$yg^*ca^* \leq yca^* \leq xa^* \leq x.$$

These two inequalities establish (4), hence (3) for E .

In applications, we will be considering left-handed Kleene algebras of monotone functions on a semilattice satisfying the ascending chain condition.

2.4 Strict Monotone Functions on a Semilattice

We model transfer functions as strict monotone functions $f : L \rightarrow L$, where L is an upper semilattice satisfying the ascending chain condition. The maps f must satisfy

$$x \leq y \Rightarrow f(x) \leq f(y) \quad (6)$$

$$f(\perp) = \perp. \quad (7)$$

There are particular strict monotone functions

$$0 = \lambda x. \perp \quad 1 = \lambda x. x.$$

The *domain* of f is the set

$$\text{dom } f = \{x \in L \mid f(x) \neq \top\}.$$

The property (6) implies that $\text{dom } f$ is closed downward under \leq .

Let K denote the family of strict monotone functions on L . We can impose a left-handed Kleene algebra structure on K as follows. First, define addition of functions pointwise:

$$\begin{aligned} (f + g)(x) &= f(x) + g(x) \\ f \leq g &\Leftrightarrow f + g = g. \end{aligned}$$

Under this definition, K forms an upper semilattice with least element 0.

Elements of K can be composed using ordinary functional composition. The operator is written \cdot and the composition of f followed by g is written fg ; thus $(fg)(x) =$

$g(f(x))$. Note that $x \in \text{dom } fg$ iff $x \in \text{dom } f$ and $f(x) \in \text{dom } g$. The identity function 1 is a two-sided identity for composition and 0 is a two-sided annihilator.

Composition distributes over $+$ on the left, but not necessarily on the right. However, it is right-subdistributive due to monotonicity. Thus K forms a left-handed idempotent semiring under the operations $+$, \cdot , 0, 1.

The element f^* is defined as the function which on input x gives the least y such that $x + f(y) \leq y$. In symbols,

$$f^*(x) = \mu y. (x + f(y) \leq y),$$

where μ is the usual least-fixpoint operator. The least fixpoint exists, since f is monotone and the ACC holds, so the monotone sequence

$$x, x + f(x), x + f(x + f(x)), \dots$$

converges after a finite number of steps, but not necessarily uniformly bounded in x ; a counterexample is given by the semilattice consisting of $\mathbb{N} \cup \{\infty\}$ with \min as join and the strict monotone function f that on input x gives ∞ if $x = \infty$, $x - 1$ if $x \geq 1$, and 0 if $x = 0$.

To show (1), we need to show that $1 + f^*f \leq f^*$, or in other words, for an arbitrary x , $x + f(f^*(x)) \leq f^*(x)$. But this is true, since $f^*(x)$ is defined to be the least element with this property.

Finally, to show (3), we need to show that $gf \leq g$ implies $gf^* \leq g$, or in other words, $f^*(g(x)) \leq g(x)$ whenever $f(g(x)) \leq g(x)$. But if $f(g(x)) \leq g(x)$, then $g(x)$ satisfies $g(x) + f(Y) \leq Y$, and $f^*(g(x))$ is the least such element.

We have shown

Theorem 2.1 *The structure $(K, +, \cdot, *, 0, 1)$ is a left-handed Kleene algebra.*

3 A Second-Order Approach

In this section we present a general second-order approach to static analysis. The technique exploits the ability to compute the Kleene algebra operations on transfer functions as defined above.

We are given a program with n instructions, and we wish to label the underlying control flow graph G of the program with elements of the semilattice L . Let E be the $n \times n$ matrix with rows and columns indexed by the vertices of G such that if (s, t) is an edge of G , then $E[s, t]$ is the transfer function labeling the edge (s, t) , and $E[s, t] = 0$ if (s, t) is not an edge of G . This matrix is easily constructed in a single pass thorough the program.

Recall from Section 2.3 that the $n \times n$ matrices over a left-handed Kleene algebra again form a left-handed Kleene algebra. We can thus speak of the matrix E^* . The entry

$E^*[u, v]$ is the join of the composition of transfer functions along all paths from u to v . The desired fixpoint dataflow labeling at any node u of G can be obtained by evaluating $E^*[s_0, u](\ell_0)$, where $\ell_0 \in L$ is the initial label of the start node s_0 . Thus the inductive labeling of the control flow graph is replaced with the computation of the matrix E^* .

A concrete implementation of this method is described in [8] in the context of Java bytecode. In that implementation, E^* is not computed directly, but rather used in conjunction with the worklist algorithm to obtain a hybrid method that uses matrix closure on a small cutset to avoid recalculation dataflow information along long paths in the control flow graph.

4 Soundness and Completeness

We argue in this section that the second-order algorithm proposed in Section 3 and the standard worklist algorithm produce the same final dataflow labeling for any type-correct program.

Let L be an upper semilattice satisfying the ACC as described in Section 2.1, and let K be the left-handed Kleene algebra of strict monotone functions $L \rightarrow L$ as described in Section 2.4. Let G be a control flow graph with vertices V , $n = |V|$, start node $s_0 \in V$, and edges labeled with transfer functions $f \in K$. Let $\ell_0 \in L$ be the initial dataflow information at s_0 .

Formally, the worklist algorithm computes a sequence of labelings $w_n : V \rightarrow L$, $n \geq 0$, as follows. We start with the initial labeling

$$w_0(u) = \begin{cases} \ell_0, & \text{if } u = s_0 \\ \perp, & \text{otherwise.} \end{cases}$$

At stage n , say we have constructed a labeling w_n . To get w_{n+1} , we take the next edge (u, v) from the worklist, apply the associated transfer function $E[u, v]$ to the current label $w_n(u)$ of u , and update the label of v with that value. Thus

$$w_{n+1}(t) = \begin{cases} E[u, v](w_n(u)) + w_n(v), & \text{if } t = v, \\ w_n(t), & \text{if } t \neq v. \end{cases}$$

The sequence w_0, w_1, \dots is monotone and converges to a fixpoint

$$w_* = \sup_n w_n.$$

This labeling is the least labeling such that for all u reachable from the start node s_0 ,

$$E[u, v](w_*(u)) \leq w_*(v) \quad (8)$$

[6]. For vertices u not reachable from s_0 , the worklist algorithm will never see them, but (8) will still hold for those

vertices, since $w_*(u) = \perp$ and $E[u, v](\perp) = \perp \leq w_*(v)$. Note that v may still be reachable from s_0 , even if u is not.

To compare this algorithm to our second-order algorithm of Section 3, it will be convenient to label vertices with certain functions $L \rightarrow L$ instead of elements of L . We lift an element $x \in L$ to an almost-constant function $\hat{x} \in K$ as follows:

$$\hat{x}(y) = \begin{cases} x, & \text{if } y \neq \perp \\ \perp, & \text{otherwise.} \end{cases}$$

Note that $\widehat{\perp} = 0$. The value x can be recovered from \hat{x} by applying \hat{x} to any element of L besides \perp . The advantage of using lifted values is that function application becomes composition, which is a Kleene algebra operation:

$$\widehat{f(x)} = \hat{x}f. \quad (9)$$

If $w : V \rightarrow L$ is a labeling, we can lift it to a second-order labeling $\widehat{w} : V \rightarrow K$ by taking

$$\widehat{w}[u] = \widehat{w(u)}.$$

For example, the lifted version of the initial labeling w_0 is

$$\widehat{w}_0[u] = \begin{cases} \widehat{\ell_0}, & \text{if } u = s_0 \\ 0, & \text{otherwise.} \end{cases}$$

Although both w and \widehat{w} are functions on V , we write $\widehat{w}[u]$ with square brackets because we will be regarding it as a row vector of length n and using it in matrix-vector computations.

We are now ready to prove our main theorem.

Theorem 4.1 *Let $w_* : V \rightarrow L$ be the final dataflow labeling produced by the worklist algorithm. Then for all $u \in V$,*

$$E^*[s_0, u](\ell_0) = w_*(u). \quad (10)$$

Proof. Let w_0 be the initial labeling, and consider the matrix-vector product $\widehat{w}_0 E^*$, which is a row vector $V \rightarrow K$. We first show that

$$\widehat{w}_0 E^* = \widehat{w}_*. \quad (11)$$

By (2), the left-hand side is the least solution of

$$\widehat{w}_0 + XE \leq X, \quad (12)$$

so it suffices to show that \widehat{w}_* is as well, since the least solution of (12) is unique. For all $u \in V$ and $x \neq \perp$,

$$\begin{aligned} \widehat{w}_0[u](x) &= w_0(u) \leq \sup_n w_n(u) \\ &= w_*(u) = \widehat{w}_*[u](x), \end{aligned}$$

therefore $\widehat{w}_0 \leq \widehat{w}_*$. Similarly,

$$\begin{aligned}
(\widehat{w}_*E)[u](x) &= \left(\sum_v \widehat{w}_*[v]E[v,u]\right)(x) \\
&= \sum_v \widehat{w}_*[v]E[v,u](x) \\
&= \sum_v E[v,u](w_*(v)) \\
&\leq w_*(u) && \text{by (8)} \\
&= \widehat{w}_*[u](x).
\end{aligned}$$

Since u and $x \neq \perp$ were arbitrary, $\widehat{w}_*E \leq \widehat{w}_*$. Since both $\widehat{w}_0 \leq \widehat{w}_*$ and $\widehat{w}_*E \leq \widehat{w}_*$, we have $\widehat{w}_0 + \widehat{w}_*E \leq \widehat{w}_*$, therefore \widehat{w}_* is a solution to (12).

To show that it is the least solution, let X be any other solution. Then $\widehat{w}_0 \leq X$. Reasoning inductively, suppose that $\widehat{w}_n \leq X$. Let (s, t) be the edge selected by the worklist algorithm at stage n . For $x \neq \perp$,

$$\begin{aligned}
\widehat{w}_{n+1}[u](x) &= w_{n+1}(u) \\
&= \begin{cases} E[v,u](w_n(v)) + w_n(u), & \text{if } (v, u) = (s, t) \\ w_n(u), & \text{otherwise} \end{cases} \\
&\leq \sum_v E[v,u](w_n(v)) + w_n(u) \\
&= \sum_v \widehat{w}_n[v]E[v,u](x) + \widehat{w}_n[u](x) \\
&= (\widehat{w}_nE)[u](x) + \widehat{w}_n[u](x) \\
&= (\widehat{w}_n(E+1))[u](x),
\end{aligned}$$

therefore $\widehat{w}_{n+1} \leq \widehat{w}_n(E+1)$. It follows that

$$\begin{aligned}
\widehat{w}_{n+1} &\leq \widehat{w}_n(E+1) \leq X(E+1) \\
&\leq XE + X \leq X.
\end{aligned}$$

Then $\widehat{w}_* = \sup_n \widehat{w}_n \leq X$, so \widehat{w}_* is the least solution of (12).

Finally, we show how (10) follows from (11). Let $u \in V$. Since $\widehat{w}_0[v] = 0$ for $v \neq s_0$, we have

$$\begin{aligned}
\widehat{w}_0[s_0]E^*[s_0, u] &= \widehat{w}_0[s_0]E^*[s_0, u] + \sum_{v \neq s_0} \widehat{w}_0[v]E^*[v, u] \\
&= \sum_v \widehat{w}_0[v]E^*[v, u] \\
&= (\widehat{w}_0E^*)[u] \\
&= \widehat{w}_*[u] && \text{by (11),}
\end{aligned}$$

thus for any $x \in L - \{\perp\}$,

$$\begin{aligned}
E^*[s_0, u](\ell_0) &= E^*[s_0, u](w_0(s_0)) \\
&= \widehat{w}_0[s_0]E^*[s_0, u](x) \\
&= \widehat{w}_*[u](x) \\
&= w_*(u). \quad \square
\end{aligned}$$

4.1 Conclusions and Future Work

We would like to implement the hybrid algorithm described in [8] and compare it experimentally to the standard worklist algorithm as specified in the Java VM specification [11]. This should not be difficult, since we already have an implementation of the latter [2].

The efficacy of our hybrid algorithm depends on finding a small cutset in the control flow graph; that is, a set of nodes intersecting every directed cycle. Finding a minimum cutset is known to be *NP*-complete, but solvable in polynomial time for reducible graphs [5]. Flowgraphs of bytecode programs compiled from Java source would ordinarily be reducible. In practice, simply taking set of all targets of back edges should give a very small cutset.

To test this, we collected some rough empirical evidence from a sample of Java bytecode programs. Of 537 programs analyzed, the median cutset size as a percentage of total program size was 2.1%. All except five programs were less than 5%. The largest program analyzed was 2668 instructions with 5 cutpoints, or 0.2%. These are very encouraging numbers indeed.

It is also apparent that our second-order method is amenable to parallelization. The worklist method is inherently sequential, since each application of a transfer function requires knowledge of its inputs, whereas compositions can be computed without knowing their inputs. This presents another intriguing possibility that we would like to investigate.

Acknowledgments

We are indebted to Stephen Chong, Andrew Myers, and Radu Rugina for valuable discussions. This work was supported in part by NSF grant CCR-0105586 and ONR Grant N00014-01-1-0968. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of these organizations or the US Government.

References

- [1] M. Abadi and R. Stata. A type system for Java bytecode subroutines. In *Proc. 25th Symp. Principles of Programming Languages*, pages 149–160. ACM SIGPLAN/SIGACT, January 1998.
- [2] F. Adelman, D. Kozen, and M. Stillerman. Malicious code detection for open firmware. In *Proc. 18th Computer Security Applications Conf. (ACSAC'02)*, pages 403–412, December 2002.
- [3] T. Ball and S. K. Rajamani. The SLAM project: Debugging system software via static analysis. In *Proc. Conf. Principles*

- of Programming Languages (POPL'02)*, pages 1–3. ACM, January 2002.
- [4] J. H. Conway. *Regular Algebra and Finite Machines*. Chapman and Hall, London, 1971.
 - [5] M. R. Garey and D. S. Johnson. *Computers and Intractibility: A Guide to the Theory of NP-Completeness*. W.H. Freeman, 1979.
 - [6] G. A. Kildall. A unified approach to global program optimization. In *Proc. Conf. Principles of Programming Languages (POPL'73)*, pages 194–206. ACM, 1973.
 - [7] S. C. Kleene. Representation of events in nerve nets and finite automata. In C. E. Shannon and J. McCarthy, editors, *Automata Studies*, pages 3–41. Princeton University Press, Princeton, N.J., 1956.
 - [8] Ł. Kot and D. Kozen. Kleene algebra and bytecode verification. Technical Report 2004-1972, Computer Science Department, Cornell University, December 2004.
 - [9] D. Kozen. A completeness theorem for Kleene algebras and the algebra of regular events. *Infor. and Comput.*, 110(2):366–390, May 1994.
 - [10] D. Kozen. Kleene algebras with tests and the static analysis of programs. Technical Report 2003-1915, Computer Science Department, Cornell University, November 2003.
 - [11] T. Lindholm and F. Yellin. *The JAVA virtual machine specification*. Addison Wesley, 1996.
 - [12] F. B. Schneider. Enforceable security policies. *ACM Trans. Information and System Security*, 3(1):30–50, February 2000.