

# Path-Sensitive Backward Slicing

Joxan Jaffar<sup>1</sup>, Vijayaraghavan Murali<sup>1</sup>, Jorge A. Navas<sup>2</sup>, and Andrew E. Santosa<sup>3</sup>

<sup>1</sup> National University of Singapore

<sup>2</sup> The University of Melbourne

<sup>3</sup> University of Sydney

**Abstract.** Backward slicers are typically path-insensitive (i.e., they ignore the evaluation of predicates in guards), or they are only partial sometimes producing too big slices. Though the value of path-sensitivity is always desirable, the major challenge is that there are, in general, an exponential number of predicate combinations to be considered. We present a *path-sensitive* backward slicer and demonstrate its practicality with real C programs. The core is a symbolic execution-based algorithm that excludes spurious dependencies lying on infeasible paths while pruning paths that cannot improve the accuracy of the dependencies already computed by other paths.

## 1 Introduction

Weiser [21] defined the *backward slice* of a program with respect to a program location  $\ell$  and a variable  $x$ , called the slicing criterion, as all statements of the program that might affect the value of  $x$  at  $\ell$ , considering all possible executions of the program. Slicing was first developed to facilitate software debugging, but it has subsequently been used for performing diverse tasks such as parallelization, software testing and maintenance, program comprehension, reverse engineering, program integration and differencing, and compiler tuning.

Although static slicing has been successfully used in many software engineering applications, slices may be quite imprecise in practice - "*slices are bigger than expected and sometimes too big to be useful [2]*". Two possible sources of imprecision are: inclusion of dependencies originated from *infeasible paths*, and merging abstract states (via join operator) along incoming edges of a *control flow merge*. A systematic way to avoid these inaccuracies is to perform path-sensitive analysis. An analysis is said to be *path-sensitive* if it keeps track of different state values based on the evaluation of the predicates at conditional branches. Although path-sensitive analyses are more precise than both flow-sensitive and context-sensitive analyses they are very rare due to the difficulty of designing efficient algorithms that can handle its combinatorial nature.

The main result of this paper is a practical path-sensitive algorithm to compute backward slices. *Symbolic execution (SE)* is the underlying technique that provides path-sensitiveness to our method. The idea behind SE is to use symbolic inputs rather than actual data and execute the program considering those symbolic inputs. During the execution of a path all its constraints are accumulated in a formula  $P$ . Whenever code of the form  $\text{if}(C) \text{ then } S1 \text{ else } S2$  is reached the execution forks the current state and updates the two copies  $P_1 \equiv P \wedge C$  and  $P_2 \equiv P \wedge \neg C$ , respectively. Then, it checks if either  $P_1$  or  $P_2$  is unsatisfiable. If yes, then the path is *infeasible* and hence, the execution stops and

backtracks to the last choice point. Otherwise, the execution continues. The set of all paths explored by symbolic execution is called the *symbolic execution tree (SET)*.

Not surprisingly, a backward slicer can be easily adapted to compute slices on SETs rather than control flow graphs (CFGs) and then mapping the results from the SET to the original CFG. It is not difficult to see that the result would be a fully path-sensitive slicer. However, there are two challenges facing this idea. First, the *path explosion problem* in path-sensitive analyses that is also present in SE since the size of the SET is exponential in the number of conditional branches. The second challenge is the infinite length of symbolic paths due to loops. To overcome the latter we borrow from [18] the use of inductive invariants produced from an abstract interpreter to automatically compute *approximate loop invariants*. Because invariants are approximate our algorithm cannot be considered fully path-sensitive in the presence of loops. Nevertheless our results in Sec. 5 demonstrate that our approach can still produce significantly more precise slices than a path-insensitive slicer.

Therefore, the main technical contribution of this paper is how to tackle the path-explosion problem. We rely on the observation that *many symbolic paths have the same impact on the slicing criterion*. In other words, there is no need to explore all possible paths to produce the most precise slice. Our method takes advantage of this observation and explores the search space by dividing the problem into smaller sub-problems which are then solved recursively. Then, it is common for many sub-problems to be “*equivalent*” to others. When this is the case, those sub-problems can be skipped and the search space can be significantly reduced with exponential speedups. In order to successfully implement this search strategy we need to (a) store the solution of a sub-problem as well as the conditions that must hold for reusing that solution, (b) reuse a stored solution if a new encountered sub-problem is “*equivalent*” to one already solved

Our approach symbolically executes the program in a depth-first search manner. This allows us to define a sub-problem as any subtree contained in the SET. Given a subtree, our method following Weiser’s algorithm computes dependencies among variables that allow us to infer which statements may affect the slicing criterion. The fundamental idea for reusing a solution is that when the set of feasible paths in a given subtree is *identical* to that of an already explored subtree, it is not possible to deduce more accurate dependencies from the given subtree. In such cases we can safely reuse dependencies from the explored subtree. However, this check is impractical because it is tantamount to actually exploring the given subtree, which defeats the purpose of reuse. Hence we define certain reusing conditions, the cornerstone of our algorithm, which are both sound and precise enough to allow reuse without exploring the given subtree.

First, we store a formula that succinctly captures all the infeasible paths detected during the symbolic execution of a subtree. We use efficient *interpolation* techniques [4] to generate *interpolants* for this purpose. Then, whenever a new subtree is encountered we check if the constraints accumulated *imply* in the logical sense the *interpolant* of an already solved subtree. If not, it means there are paths in the new subtree which were unexplored (infeasible) before, and so we need to explore the subtree in order to be sound. Otherwise, the set of paths in the new subtree is a *subset* of that of the explored subtree. However, being a subset is not sufficient for reuse since we need to know if they are *equivalent*, but the equivalence test, as mentioned before, is impractical.

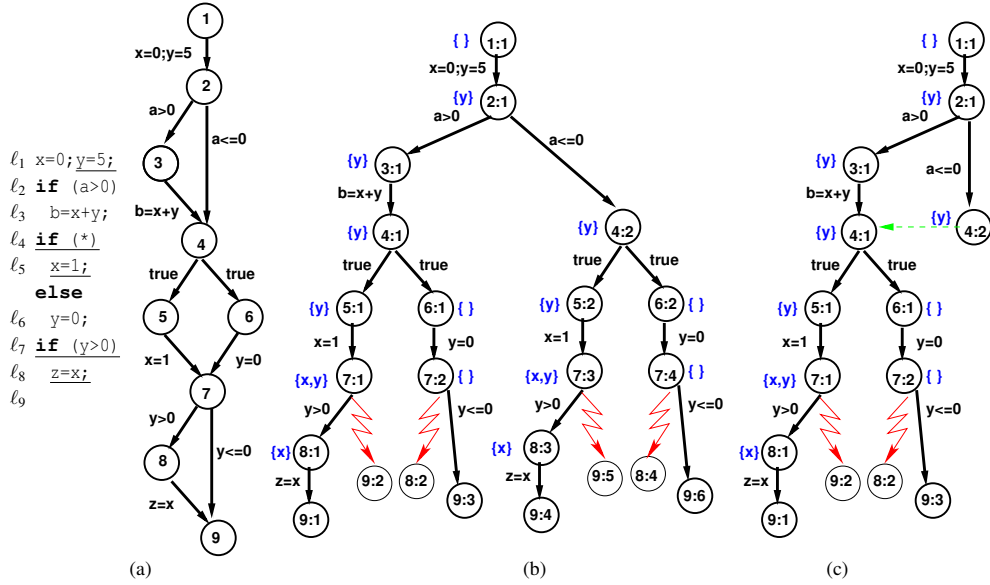


Fig. 1. (a) A program and its transition system, (b) its naive symbolic execution tree (SET) and (c) its interpolation-based SET, for slicing criterion (underlined statements)  $\langle \ell_9, \{z\} \rangle$

Here, we make use of our intuition that only few paths contribute to the dependency information in every subtree. Hence, to check for equivalence of subtrees we need not check all paths, but only those that contribute to the dependencies, what we call the *witness paths*. Now, if the implication succeeds we also check if the witness paths of the explored subtree are feasible in the new subtree. If yes, we reuse dependencies. Otherwise, the equivalence test failed.

Finally, as we will discuss in Sec. 6, some previous works have tackled the problem of path-sensitive backward slicing before. However, to the best of our knowledge either they suffer from the path-explosion problem or scalability is achieved at the expense of losing some path-sensitiveness. One essential result of our method is that it produces *exact* slices for loop-free programs. By “exact” we mean that the algorithm guarantees to not produce dependencies from spurious<sup>1</sup> (i.e., non-executable) paths. In other words, it produces the *smallest* possible, *sound* slice of a loop-free program for any given slicing criterion. Our method mitigates the path-explosion problem using a combination of interpolants and witness paths that allows pruning significantly the search space.

## 2 Motivating Example

We first describe our approach through an example. Consider the program in Fig. 1(a) and assume we would like to slice it wrt location  $\ell_9$  and variable  $z$ . The statement  $x:=0$

<sup>1</sup> Of course, limited by theorem prover technology which decides whether a formula is unsatisfiable or not.

at  $\ell_1$  should not be included in the slice because any path that reaches  $\ell_8$  through  $\ell_5$  redefines  $x$  and any path that reaches  $\ell_8$  through  $\ell_6$  (without redefining  $x$ ) is infeasible. Note that a path insensitive algorithm would not be able to infer this from the CFG.

Fig. 1(b) shows the naive symbolic execution tree of the program. The nodes are labeled with  $\ell : k$  ( $\ell$  is a program location and  $k$  is an identifier to distinguish nodes with the same program location belonging to different symbolic paths) and edges between two locations are labeled by the intervening program operation. Black (solid) edges denote feasible transitions and red (zigzag) edges denote infeasible transitions. Each node is annotated with its *dependency set* in blue (between brackets) obtained by running, for instance, Weiser’s [21] algorithm. Informally, a dependency set at location  $\ell$  contains all variables that may affect the slicing criterion from any path reachable from  $\ell$ . A statement at  $\ell$  is included in the slice if the intersection between the dependency set at  $\ell$  and the set of variables defined at  $\ell$  (i.e., left-hand side of the assignment) is not empty. Note that the dependency set at 2:1 only contains  $y$  and therefore, the statement  $x:=0$  at  $\ell_1$  would not be included in the slice. Hence it is clear that the path-sensitive SET improves the accuracy of slices. The problem is that the size of the tree is exponential in the number of branches. However, consider now the tree in Fig. 1(c) constructed by our method<sup>2</sup> where green (dotted) edges denote reusing transitions. This tree contains the same relevant information needed to exclude  $x:=0$  from the slice but without some redundant paths present in Fig. 1(b), achieving exponential savings. Let us see how the tree in Fig. 1(c) is generated.

Our algorithm performs symbolic execution guided by depth-first search exploring first the path  $\pi \equiv \ell_1 \cdot \ell_2 \cdot \ell_3 \cdot \ell_4 \cdot \ell_5 \cdot \ell_7 \cdot \ell_8 \cdot \ell_9$ . As usual, it accumulates the constraints along the path in a formula  $\Pi$ , where variable redefinitions are denoted by primed versions. For the above path,  $\Pi_{9,1} \equiv x = 0 \wedge y = 5 \wedge a > 0 \wedge b = x + y \wedge x' = 1 \wedge y > 0 \wedge z = x'$  is the formula built at 9:1, which is satisfiable. It then applies Weiser’s algorithm to compute the dependency set at each node along the path. In addition, it also computes at each node one of the reusing conditions: the (smallest possible) set of paths from which the dependency set was generated. For example, at 7:1 the dependency set  $\{x, y\}$  was obtained from the path  $\ell_7 \cdot \ell_8 \cdot \ell_9$ , at 4:1 the dependency set  $\{y\}$  was obtained from  $\ell_4 \cdot \ell_5 \cdot \ell_7 \cdot \ell_8 \cdot \ell_9$ , and so on. These paths are called the *witness paths* and they represent the paths along which each variable in the dependency set affects the slicing criterion.

Next our algorithm backtracks and explores the path  $\pi \equiv \ell_1 \cdot \ell_2 \cdot \ell_3 \cdot \ell_4 \cdot \ell_5 \cdot \ell_7 \cdot \ell_9$  with constraints  $\Pi_{9,2} \equiv x = 0 \wedge y = 5 \wedge a > 0 \wedge b = x + y \wedge x' = 1 \wedge y \leq 0$ . This formula is unsatisfiable and hence the path is infeasible. Now it generates another reusing condition: a formula called the *interpolant* that captures the essence of the reason of infeasibility of the path. The main purpose of the interpolant is to exclude irrelevant facts pertaining to the infeasibility so that the reusing conditions are more likely to be reusable in future. For the above path a possible interpolant is  $y = 5$  which is enough to capture its infeasibility and the infeasibility of any path that carries the constraint  $y \leq 0$ . In summary, our algorithm generates two reusing conditions: witness paths from feasible paths and interpolants from infeasible paths.

Next it backtracks and explores the path  $\pi \equiv \ell_1 \cdot \ell_2 \cdot \ell_3 \cdot \ell_4 \cdot \ell_6 \cdot \ell_7$ . At 7:2, it checks whether it can reuse the solution from 7:1 by checking if the accumulated constraints

<sup>2</sup> In fact, it is a Directed Acyclic Graph (DAG) due to the existence of reusing edges.

$\Pi_{7:2} \equiv x = 0 \wedge y = 5 \wedge a > 0 \wedge b = x + y \wedge y' = 0$  imply the interpolant at 7:1,  $y' = 5^3$ . Since the implication fails, it has to explore 7:2 in order to be sound. The subtree after exploring this can be seen in Fig. 1(c). An important thing to note here is that while applying Weiser’s algorithm, it has obtained a more accurate dependency set (empty set) at 7:2 than that which would have been obtained if it reused the solution from 7:1. Also note that at 4:1, the dependency set is still  $\{y\}$  with witness path  $\ell_4 \cdot \ell_5 \cdot \ell_7 \cdot \ell_8 \cdot \ell_9$  and interpolant  $y = 5$ .

Note what happens now. When our algorithm backtracks to explore the path  $\pi \equiv \ell_1 \cdot \ell_2 \cdot \ell_4$ , it checks at 4:2 if it can reuse the solution from 4:1. This time, the accumulated constraints  $x = 0 \wedge y = 5 \wedge a \leq 0$  imply the interpolant at 4:1,  $y = 5$ . In addition, the witness path at 4:1 is also feasible under 4:2. Hence, it simply reuses the dependency set  $\{y\}$  from 4:1 both in a *sound* and *precise* manner, and backtracks without exploring 4:2. In this way, it achieves exponential savings while still maintaining as much as accuracy as the naive SET in Fig. 1(b). Now, when Weiser’s algorithm propagates back the dependency set  $\{y\}$  from 4:2, we get the dependency set  $\{y\}$  again at 2:1, and the statement  $x:=0$  at 1:1 is not included in the slice.

### 3 Background

**Syntax.** We restrict our presentation to a simple imperative programming language where all basic operations are either assignments or assume operations, and the domain of all variables are integers. The set of all program variables is denoted by *Vars*. An *assignment*  $x := e$  corresponds to assign the evaluation of the expression  $e$  to the variable  $x$ . In the *assume* operator,  $\text{assume}(c)$ , if the boolean expression  $c$  evaluates to *true*, then the program continues, otherwise it halts. The set of operations is denoted by *Ops*. We then model a program by a *transition system*. A transition system is a quadruple  $\langle \Sigma, I, \longrightarrow, O \rangle$  where  $\Sigma$  is the set of states and  $I \subseteq \Sigma$  is the set of initial states.  $\longrightarrow \subseteq \Sigma \times \Sigma \times Ops$  is the transition relation that relates a state to its (possible) successors executing operations. This transition relation models the operations that are executed when control flows from one program location to another. We shall use  $\ell \xrightarrow{\text{op}} \ell'$  to denote a transition relation from  $\ell \in \Sigma$  to  $\ell' \in \Sigma$  executing the operation  $\text{op} \in Ops$ . Finally,  $O \subseteq \Sigma$  is the set of final states.

**Symbolic Execution.** A *symbolic state*  $v$  is a triple  $\langle \ell, s, \Pi \rangle$ . The symbol  $\ell \in \Sigma$  corresponds to the current program location (with special symbols for initial location,  $\ell_{\text{start}}$ , and final location,  $\ell_{\text{end}}$ ). The symbolic store  $s$  is a function from program variables to terms over input symbolic variables. Each program variable is initialized to a fresh input symbolic variable. The *evaluation*  $\llbracket e \rrbracket_s$  of an arithmetic expression  $e$  in a store  $s$  is defined as usual:  $\llbracket v \rrbracket_s = s(v)$ ,  $\llbracket n \rrbracket_s = n$ ,  $\llbracket e + e' \rrbracket_s = \llbracket e \rrbracket_s + \llbracket e' \rrbracket_s$ ,  $\llbracket e - e' \rrbracket_s = \llbracket e \rrbracket_s - \llbracket e' \rrbracket_s$ , etc. The evaluation of Boolean expression  $\llbracket b \rrbracket_s$  can be defined analogously. Finally,  $\Pi$  is called *path condition* and it is a first-order formula over the symbolic inputs and it accumulates constraints which the inputs must satisfy in order for an execution to follow the particular corresponding path. The set of first-order formulas and symbolic states are denoted by *FOL* and *SymStates*, respectively. Given a transition system  $\langle \Sigma, I, \longrightarrow, O \rangle$

<sup>3</sup> The interpolant always considers the latest versions of the variables.

and a state  $\mathfrak{v} \equiv \langle \ell, s, \Pi \rangle \in \text{SymStates}$ , the symbolic execution of  $\ell \xrightarrow{\text{op}} \ell'$  returns another symbolic state  $\mathfrak{v}'$  defined as:

$$\mathfrak{v}' \triangleq \begin{cases} \langle \ell', s, \Pi \wedge \llbracket c \rrbracket_s \rangle & \text{if } \text{op} \equiv \text{assume}(c) \text{ and } \Pi \wedge \llbracket c \rrbracket_s \text{ is satisfiable} \\ \langle \ell', s[x \mapsto \llbracket e \rrbracket_s], \Pi \rangle & \text{if } \text{op} \equiv x := e \end{cases} \quad (1)$$

Note that Eq. (1) queries a *theorem prover* for satisfiability checking on the path condition. We assume the theorem prover is sound but not necessarily complete. That is, the theorem prover must say a formula is unsatisfiable only if it is indeed so.

Abusing notation, given a symbolic state  $\mathfrak{v} \equiv \langle \ell, s, \Pi \rangle$  we define  $\llbracket \mathfrak{v} \rrbracket : \text{SymStates} \rightarrow \text{FOL}$  as the formula  $(\bigwedge_{v \in \text{Vars}} \llbracket v \rrbracket_s) \wedge \Pi$  where  $\text{Vars}$  is the set of program variables.

A *symbolic path*  $\pi \equiv \mathfrak{v}_0 \cdot \mathfrak{v}_1 \cdot \dots \cdot \mathfrak{v}_n$  is a sequence of symbolic states such that  $\forall i \bullet 1 \leq i \leq n$  the state  $\mathfrak{v}_i$  is a *successor* of  $\mathfrak{v}_{i-1}$ . A symbolic state  $\mathfrak{v}' \equiv \langle \ell', \cdot, \cdot \rangle$  is a successor of another  $\mathfrak{v} \equiv \langle \ell, \cdot, \cdot \rangle$  if there exists a transition relation  $\ell \xrightarrow{\text{op}} \ell'$ . A path  $\pi \equiv \mathfrak{v}_0 \cdot \mathfrak{v}_1 \cdot \dots \cdot \mathfrak{v}_n$  is *feasible* if  $\mathfrak{v}_n \equiv \langle \ell, s, \Pi \rangle$  such that  $\llbracket \Pi \rrbracket_s$  is satisfiable. If  $\ell \in O$  and  $\mathfrak{v}_n$  is feasible then  $\mathfrak{v}_n$  is called *terminal* state. Otherwise, if  $\llbracket \Pi \rrbracket_s$  is unsatisfiable the path is called *infeasible* and  $\mathfrak{v}_n$  is called an *infeasible* state. If there exists a feasible path  $\pi \equiv \mathfrak{v}_0 \cdot \mathfrak{v}_1 \cdot \dots \cdot \mathfrak{v}_n$  then we say  $\mathfrak{v}_k$  ( $0 \leq k \leq n$ ) is *reachable* from  $\mathfrak{v}_0$  in  $k$  steps. We say  $\mathfrak{v}''$  is reachable from  $\mathfrak{v}$  if it is reachable from  $\mathfrak{v}$  in some number of steps.

A *symbolic execution tree* contains all the execution paths explored during the symbolic execution of a transition system by triggering Eq. (1). The nodes represent symbolic states and the arcs represent transitions between states.

**Program Slicing via Abstract Interpretation.** The *backward slice* of a program wrt a program location  $\ell$  and a set of variables  $V \subseteq \text{Vars}$ , called the *slicing criterion*  $\langle \ell, V \rangle$ , is all statements of the program that might affect the values of  $V$  at  $\ell$ .<sup>4</sup> We follow the dataflow approach described by Weiser [21] reformulated as an abstract domain  $\mathcal{D} \equiv \{\perp\} \cup \mathcal{P}(\text{Vars})$  (where  $\mathcal{P}(\text{Vars})$  is the powerset of program variables) with a lattice structure  $\langle \sqsubseteq, \perp, \sqcup, \sqcap, \top \rangle$ , such that  $\sqsubseteq \equiv \subseteq$ ,  $\sqcup \equiv \cup$ , and  $\sqcap \equiv \cap$  are conveniently lifted to consider the element  $\perp$ .

We say  $\sigma_\ell \in \mathcal{D}$  is the approximate set of variables at location  $\ell$  that may affect the slicing criterion. We will abuse notation to denote the dependencies associated to a symbolic state  $\mathfrak{v}$  also as  $\sigma_{\mathfrak{v}}$ . *Backward data dependencies* can be formulated using this set, defining two kinds of dataflow information. Given a transition relation  $\ell \xrightarrow{\text{op}} \ell'$  we define  $\text{def}(\text{op})$  and  $\text{use}(\text{op})$  as the sets of variables altered and used during the execution of  $\text{op}$ , respectively. Then,

$$\sigma_\ell \triangleq \begin{cases} (\sigma_{\ell'} \setminus \text{def}(\text{op})) \cup \text{use}(\text{op}) & \text{if } \sigma_{\ell'} \cap \text{def}(\text{op}) \neq \emptyset \\ \sigma_{\ell'} & \text{otherwise} \end{cases} \quad (2)$$

where  $\sigma_{\ell'} = V$  if  $\ell' = \ell_{\text{end}}$ . We say a transition relation  $\ell \xrightarrow{\text{op}} \ell'$  where  $\text{op} \equiv x := e$  is included in the slice if:

$$\sigma_{\ell'} \cap \text{def}(\text{op}) \neq \emptyset \quad (3)$$

*Backward control dependencies* can also affect the slicing criterion. A transition relation  $\delta \equiv \ell \xrightarrow{\text{op}} \ell'$  where  $\text{op} \equiv \text{assume}(c)$  is included in the slice if any transition relation

<sup>4</sup> W.l.o.g., we assume in this paper a single slicing criterion at  $\ell_{\text{end}}$ .

under the range of influence<sup>5</sup> (the function INFL will compute the range of influence) of  $\delta$  is included in the slice, and

$$\sigma_\ell \triangleq \sigma_{\ell'} \cup use(op) \quad (5)$$

Finally, a function  $\widehat{pre}_{\mathcal{D}}(\sigma_\ell, op)$  that returns the *pre-state* after executing backwards the operation  $op$  with the *post-state*  $\sigma_\ell$  is defined using Eqs. (2,3,4,5).

## 4 Algorithm

A path-sensitive slicing algorithm over a symbolic execution tree (SET) can be defined as an *annotation* process which labels each symbolic state  $\nu \equiv \langle \ell, \cdot, \cdot \rangle$  with  $\sigma_\ell \in \mathcal{D}$  by computing a fixpoint (later formalized) over the tree, using Eqs. (2,5) described in Sec. 3. In an interleaved process, the final SET is obtained through Eqs. (3,4). Since the SET may have multiple instances of the same transition relation, we say that a transition relation is included in the final slice if at least one of its instances is included in the slice on the SET. It is easy to see that the path-sensitiveness comes from how symbolic execution builds the tree since no dependencies from a non-executable path can be considered.

Our algorithm performs symbolic execution in a depth-first search manner excluding all infeasible paths. Whenever the forward traversal of a path finishes due to a (a) terminal state, (b) infeasible state, or (c) reusing state (i.e., a state reusing a solution from another state), the algorithm halts and backtracks to the next path. During this backtracking each symbolic state  $\nu$  is labelled with its *solution*, i.e., the set of variables  $\sigma_\nu$  at  $\nu$  that may affect the slicing criterion. Furthermore, the reusing conditions are computed at each state for future use. We first introduce formally the two key concepts which will decide whether a solution can be reused or not.

**Definition 1 (Interpolant).** *Given a pair of first order logic (FOL) formulas  $A$  and  $B$  such that  $A \wedge B$  is false a Craig interpolant [4] wrt  $A$  is another FOL formula  $\overline{\Psi}$  such that (a)  $A \models \overline{\Psi}$ , (b)  $\overline{\Psi} \wedge B$  is false, and (c)  $\overline{\Psi}$  is formed using common variables of  $A$  and  $B$ .*

Note that interpolation allows us to remove irrelevant facts from  $A$  without affecting the unsatisfiability of  $A \wedge B$ .

**Definition 2 (Witness Paths and Formulas).** *Given a symbolic state  $\nu \equiv \langle \ell, \cdot, \cdot \rangle$  annotated with the set of variables  $\sigma_\nu$  that affect the slicing criterion at  $\ell_{end}$ , a witness path for a variable  $v \in \sigma_\nu$  is a symbolic path  $\pi \equiv \langle \ell, \cdot, \cdot \rangle \dots \langle \ell_{end}, \cdot, \Pi_{end} \rangle$  with the final symbolic state  $\nu' \equiv \langle \ell_{end}, \cdot, \Pi_{end} \rangle$  such that  $\llbracket \nu' \rrbracket$  is satisfiable (i.e.,  $\pi$  is feasible). We call  $\llbracket \nu' \rrbracket$  the witness formula of  $v$ , denoted  $\omega_v$ .*

Intuitively, a witness path for a variable at a node is a path below the node along which the variable affects the slicing criterion at the end. A witness formula represents a condition sufficient for the variable to affect the slicing criterion along the witness path.

Prior to establishing the reusing conditions, we augment the abstract domain  $\mathcal{D}$  to accommodate the witness formulas. Here, and in the rest of the paper, we will refer to the term “dependency” as the set of variables that may affect the slicing criterion together with their witnesses.

<sup>5</sup> More formally, the range of influence for  $\delta$  is the set of transition relations defined in any path from  $\delta$  to its *nearest postdominator* in the transition system.

$$- \sqcup : \mathcal{D}^\omega \times \mathcal{D}^\omega \rightarrow \mathcal{D}^\omega$$

$$\sigma^{\omega_1} \sqcup \sigma^{\omega_2} \triangleq \sigma^{\omega_1} \cup \sigma^{\omega_2}$$

$$- \sqsubseteq : \mathcal{D}^\omega \times \mathcal{D}^\omega \rightarrow \text{Bool}$$

$$\sigma^{\omega_1} \sqsubseteq \sigma^{\omega_2} \text{ if and only if } \sigma^{\omega_1} \subseteq \sigma^{\omega_2}$$

$$- \widehat{pre} : \mathcal{D}^\omega \times (\Sigma \times \Sigma \times \text{Ops}) \times (\text{Vars} \rightarrow \text{SymVars}) \rightarrow \mathcal{D}^\omega.$$

$$\widehat{pre}(\sigma^{\omega'}, \ell \xrightarrow{\text{op}} \ell', s) \triangleq \left\{ \begin{array}{l} \text{let } \sigma^\omega := \widehat{pre\_aux}(\sigma^{\omega'}, \ell \xrightarrow{\text{op}} \ell', s) \\ \text{foreach } \langle x, \omega_x \rangle \in \sigma^\omega, \langle x, \omega_{x'} \rangle \in \sigma^{\omega'} \\ \sigma^\omega := \sigma^\omega \setminus \{ \langle x, \omega_x \rangle, \langle x, \omega_{x'} \rangle \} \\ \text{if } \omega_x \models \omega_{x'} \text{ then } \sigma^\omega := \sigma^\omega \cup \{ \langle x, \omega_{x'} \rangle \} \\ \text{else } \sigma^\omega := \sigma^\omega \cup \{ \langle x, \omega_x \rangle \} \\ \text{if } (\sigma^\omega \cap \text{def}(\text{op}) \text{ or } \text{INFL}(\ell \rightarrow \ell') \cap \mathcal{S} \neq \emptyset) \text{ then} \\ \mathcal{S} := \mathcal{S} \cup \{ \ell \rightarrow \ell' \} \\ \text{in } \sigma^\omega \end{array} \right.$$

where:

$$\widehat{pre\_aux}(\sigma^{\omega'}, \ell \xrightarrow{\text{op}} \ell', s) \triangleq \left\{ \begin{array}{l} \{ \langle x, \omega_x \wedge \llbracket y = e \rrbracket_s \rangle \mid \langle x, \omega_x \rangle \in \sigma^{\omega'}, \text{op} \equiv y := e, x \notin \text{def}(\text{op}) \} \cup \\ \{ \langle v, \omega_x \wedge \llbracket y = e \rrbracket_s \rangle \mid \langle x, \omega_x \rangle \in \sigma^{\omega'}, \text{op} \equiv y := e, x \in \text{def}(\text{op}), v \in \text{use}(\text{op}) \} \cup \\ \{ \langle x, \omega_x \wedge \llbracket c \rrbracket_s \rangle \mid \langle x, \omega_x \rangle \in \sigma^{\omega'}, \text{op} \equiv \text{assume}(c) \} \cup \\ \{ \langle x, \llbracket \Pi \pi \rrbracket_s \wedge \llbracket c \rrbracket_s \rangle \mid \langle x, \cdot \rangle \notin \sigma^{\omega'}, \text{op} \equiv \text{assume}(c), x \in \text{use}(\text{op}), \\ \text{INFL}(\ell \rightarrow \ell') \cap \mathcal{S} \neq \emptyset, \exists \pi \equiv \ell' \cdot \dots \cdot \ell_{end} \} \end{array} \right.$$

**Fig. 2.** Main Abstract Operations for  $\mathcal{D}^\omega$

**Definition 3** ( $\mathcal{D}^\omega$ ). We define a new abstract domain  $\mathcal{D}^\omega$  as a lattice  $\langle \sqsubseteq, \perp, \sqcup, \top \rangle$  such that  $\mathcal{D}^\omega \triangleq \{ \perp \} \cup \mathcal{P}(\text{Vars} \times \text{FOL})$  (i.e., set of pairs of the form  $\langle x, \omega_x \rangle$  where  $x$  is a variable and  $\omega_x$  is its witness formula) and abstract operations described in Fig. 2.<sup>6</sup>

Note that the witness formulas can be obtained only from (feasible) paths in the program. Therefore, the number of witness formulas is always finite. As we will see later, even with loops, the size of each witness formula is also finite because we make the symbolic subtree of the loop finite. That is, we perform symbolic execution on a finite program once loop invariants are given. This ensures that the abstract domain  $\mathcal{D}^\omega$  is finite and hence, termination is guaranteed for any fixpoint computation based on it.

In Fig. 2, the operator  $\sqcup$  computes the least upper bound of the abstract states by simply applying the set union of the two set of states. The operator  $\sqsubseteq$  simply tests whether one set is a subset of the other.  $\widehat{pre}$  is a bit more elaborated but basically consists of the Eqs. (2,3,4,5) defined in Sec. 3 extended with witnesses formulas. We assume here and in the algorithm in Fig. 3 that  $\widehat{pre}$  accesses  $\mathcal{S}$  which is the set of transitions included in the slice. In function  $\widehat{pre\_aux}$ , there are four cases to handle different kinds of statements and dependencies:

- In the first two cases, if the operation is an assignment, the dependencies are propagated from the *defined* to the *used* variables and any dependency from a variable not *defined* is kept. In these cases, the pre-state witness formula is the conjunction of the post-state witness formula with the corresponding statement.

<sup>6</sup> For lack of space, trivial treatment of the element  $\perp$  is omitted from operations in Fig. 2.



- In the third case, if the operation is an assume, any *used* variable is preserved, with its pre-state witness formula being the conjunction of the post-state witness formula and the corresponding guard.
- In the last case, for any variable  $x$  occurring in an assume statement without any dependency, if any transition under the range of influence (computed by INFL) of the assume is already in the slice, then  $x$  is added (due to control dependency) and its witness formula is the conjunction of the guard and the path condition of any (feasible) path from the assume statement that leads to the end of the program.

In addition, in function  $\widehat{pre}$  whenever two pairs from the set of dependencies computed by  $\widehat{pre\_aux}$  refer to the same variable, we choose the one with the weaker witness formula (which is more likely to be reused). Finally, a transition is included in the slice if one of the Eqs. (3,4) holds.

**Definition 4 (Reusing Conditions).** *Given a current symbolic state  $\nu \equiv \langle \ell, \cdot, \Pi \rangle$  and an already solved symbolic state  $\nu' \equiv \langle \ell', \cdot, \cdot \rangle$  such that  $\overline{\Psi}$  is the interpolant generated for  $\nu'$  and  $\sigma^\omega$  are the dependencies together with their attached witnesses at  $\nu'$ , we say  $\nu$  is equivalent to  $\nu'$  (or  $\nu$  can reuse the solution at  $\nu'$ ) if the following conditions hold:*

(a)  $\llbracket \nu \rrbracket \models \overline{\Psi}$     (b)  $\forall \langle x, \cdot \rangle \in \sigma^\omega \bullet \exists \langle x, \omega_x \rangle \in \sigma^\omega$  such that  $\llbracket \nu \rrbracket \wedge \omega_x$  is satisfiable (6)

The condition (a) affects *soundness* and it ensures that the set of symbolic paths reachable from  $\nu$  must be a subset of those from  $\nu'$ . The condition (b) is the witness check which essentially states that for each variable  $x$  in the dependency set at  $\nu'$ , there must be at least one witness path with formula  $\omega_x$  that is feasible from  $\nu$ . This affects *accuracy* and ensures that the reuse of dependencies does not incur any loss of precision.

We now describe in detail the main features of our algorithm defined by the function  $\text{BackwardDeps}_V$  in Fig. 3. The main purpose of  $\text{BackwardDeps}_V$  is to keep track of the *backward dependencies* between the program variables and the slicing criterion by inferring for each state the set of variables that may affect the slicing criterion. From these dependencies it is straightforward to obtain the slice of the program as explained at the beginning of this section. For clarity of presentation, let us omit the content of the *grey* boxes and assume programs do not have loops, which we will come to later.

$\text{BackwardDeps}_V : \text{SymStates} \times \mathcal{D}^\omega \rightarrow \text{FOL} \times \mathcal{D}^\omega \times \text{Bool}$  requires the program to have been translated to a transition system  $\langle \Sigma, I, \longrightarrow, O \rangle$  and takes as input an initial symbolic state  $\nu \equiv \langle \ell \in I, \varepsilon, \text{true} \rangle$  and an initially empty  $\sigma^\omega$ .  $V$  is the set of variables of the slicing criterion. The set of transitions included in the slice,  $\mathcal{S}$ , is also empty. Recall that  $\mathcal{S}$  is only modified by  $\widehat{pre}_{\mathcal{D}^\omega}$ , and hence, we omit it from the description of the algorithm defining it as a global variable. The output is a triple with the interpolant, dependencies (i.e., reusing conditions and solution) and a boolean flag representing whether any change occurred in a dependency set at any symbolic state during the algorithm's backward traversal (this is used mainly to handle loops later). The actual object of interest computed by the algorithm is the set of transitions  $\mathcal{S}$  included in the slice.

$\text{BackwardDeps}_V$  implements a recursive algorithm whose objective is to generate a finite complete SET while reusing solutions whenever possible to avoid path explosion. Line 1 initializes the (local) variable *change* to *false*, which will be updated later. Next, the three base cases for symbolic states are handled - infeasible, terminal, and reuse:

```

BackwardDepsV(v ≡ ⟨ℓ, s, Π⟩, σω)
1: change := false
2: if INFEASIBLE(v) then ⟨ $\bar{\Psi}$ , σω⟩ := ⟨false, ∅⟩ and goto 13
3: if TERMINAL(v) then ⟨ $\bar{\Psi}$ , σω⟩ := ⟨true, {⟨v, true⟩ | v ∈ V}⟩ and goto 13
4: if ∃ v' ≡ ⟨ℓ, s, ·⟩ labelled with ⟨ $\bar{\Psi}$ , σω⟩ such that REUSE(v, v') then goto 13
5: if ℓ is the header of a loop then
6:    $\bar{v}$  := invariant(v, ℓ → ... → ℓ)
7:   ⟨ $\bar{\Psi}$ , σω, change⟩ := UnwindTreeV( $\bar{v}$ , σω) and goto 13
8: if ∃ ℓ' such that ℓ → ℓ' is a backedge of a loop then
9:   ⟨·, ·,  $\bar{\Pi}$ ⟩ := invariant(v, ℓ' → ... → ℓ)
10:  ⟨ $\bar{\Psi}$ , σω⟩ := ⟨ $\bar{\Pi}$ , σω⟩ and goto 13
11: ⟨ $\bar{\Psi}$ , σω, change⟩ := UnwindTreeV(v, σω)
12: change := change ∨ v is labeled with ⟨·, σωold⟩⟩ such that ¬(σωold ⊆Dω σω)
13: label v with ⟨ $\bar{\Psi}$ , σω⟩ and return ⟨ $\bar{\Psi}$ , σω, change⟩

```

```

UnwindTreeV(v ≡ ⟨ℓ, s, Π⟩, σωin)
1:  $\bar{\Psi}$  := true, σω := σωin, change := false
2: foreach transition relation ℓ  $\xrightarrow{op}$  ℓ'
3:   v' ≜ { ⟨ℓ', s, Π ∧ [[c]]s⟩           if op ≡ assume(c)
          ⟨ℓ', s[x ↦ Sx], Π ∧ [[x = e]]s⟩ if op ≡ x := e and Sx fresh variable
4:   ⟨ $\bar{\Psi}'$ , σω', c⟩ := BackwardDepsV(v', σωin)
5:    $\bar{\Psi}$  :=  $\bar{\Psi}$  ∧ wlp(op,  $\bar{\Psi}'$ )
6:   σω := σω ⊔Dω preDω(σω', op, s)
7:   change := change ∨ c
8: return ⟨ $\bar{\Psi}$ , σω, change⟩

```

```

BackwardDepsLoopV(v, σω)
1: σω' := σω, change := false
2: do ⟨·, σω', change⟩ := BackwardDepsV(v, σω') while change end

```

**Fig. 3.** Path-Sensitive Backward Slicing Analysis

- In line 2, the function INFEASIBLE(⟨·, ·, Π⟩) checks whether Π is satisfiable. If not, the symbolic execution detects an infeasible path and halts, excluding any dependency which would have been inferred from the non-executable path. In addition, it produces an interpolant from Π and *false*, namely  $\bar{\Psi} \equiv false$ , which generalizes the current path condition (Π ⊨  $\bar{\Psi}$  and  $\bar{\Psi}$  is *false*). Since the path is not executable there is no variable that may affect the slicing criterion and hence, the set of dependencies returned is empty.
- In line 3, the function TERMINAL(⟨ℓ, ·, ·⟩) checks if the symbolic state is a terminal node by checking if ℓ = ℓ<sub>end</sub>. If yes, the execution has reached the end of a path. Since the path is feasible, it can be fully generalized returning the interpolant  $\bar{\Psi} \equiv true$ . Since ℓ is a terminal node, the set of dependencies is the set of variables in the slicing criterion, V. The witness formula for each variable from V is initially *true*.

- In line 4 the algorithm searches for another state  $v'$  whose dependencies can be reused by the current state  $v$  so that the symbolic execution can be stopped. For this, the function  $\text{REUSE}(v, v')$  tests both the reusing conditions in Eq. 6. If the test holds, the state  $v$  can reuse the dependencies computed by  $v'$ .

If all three base cases fail, the algorithm unwinds the execution tree by calling the procedure  $\text{UnwindTree}_V$  at line 11.  $\text{UnwindTree}_V$ , at line 3, executes one symbolic step<sup>7</sup> and calls the main procedure  $\text{BackwardDeps}_V$  with the successor state (line 4). After the call the two key remaining steps are to compute:

- the interpolant  $\bar{\Psi}$  ( $\text{UnwindTree}_V$  line 5) that generalizes the symbolic execution tree below  $v$  while preserving its infeasible paths. The procedure  $\widehat{wlp} : \text{Ops} \times \text{FOL} \rightarrow \text{FOL}$  ideally computes the *weakest liberal precondition* ( $wlp$ ) [7] which is the weakest formula on the initial state ensuring the execution of  $op$  results in a final state  $\bar{\Psi}$ . In practice, we approximate  $wlp$  by making a linear number of calls to a theorem prover following techniques described in [14]. The interpolant  $\bar{\Psi}$  is an FOL formula consisting of the conjunction of the result of  $\widehat{wlp}$  on each child's interpolant.
- the solution,  $\sigma^0$ , for the current state  $v$  at line 6 which is computed by executing  $\widehat{pre}_{\mathcal{D}^\omega}$  on each child's solution and then combining all solutions using  $\sqcup_{\mathcal{D}^\omega}$ .

In addition, at line 7 it also records changes in any child's symbolic state (if any) and then returns a triple in the same format as  $\text{BackwardDeps}_V$ 's return value. In  $\text{BackwardDeps}_V$ , line 12 updates  $\text{change}$  to *true* if either it was set to *true* in  $\text{UnwindTree}_V$  at line 11 or the current symbolic state is about to be updated with a more precise solution than that it already has. The final operation before returning from  $\text{BackwardDeps}_V$  is to label the state  $v$  with the reusing conditions and solution (line 13).

Now we continue describing our algorithm by discussing how it handles loops. The main issue is to produce a finite symbolic execution tree on which a fixpoint of the dependencies can be computed.

For this, the algorithm in Fig. 3 takes an annotated transition system in which program points are labelled with inductive invariants inferred automatically by an abstract interpreter using an abstract domain such as *octagons* or *polyhedra* (we borrow the ideas presented in [18] for this purpose). We assume the abstract interpreter provides a function  $\text{getAssrt}$  which, given a program location  $\ell$  and a symbolic store  $s$ , returns an assertion in the form of an FOL formula renamed using  $s$ , which holds at  $\ell$ . Note that when applied at loop headers,  $\text{getAssrt}$  will return a loop invariant. However, we would like to strengthen it using the constraints propagated from the symbolic execution. The function  $\text{invariant}$  performs this task as follows:

$$\text{invariant}(\langle \ell, s, \Pi \rangle, \ell_1 \rightarrow \ell_n) \triangleq \begin{cases} \text{let } s' := \text{havoc}(s, \text{modifies}(\ell_1 \rightarrow \ell_n)) \\ \quad \bar{\Pi} := \text{getAssrt}(\ell, s') \wedge \Pi \\ \text{in } \langle \ell, s', \bar{\Pi} \rangle \end{cases}$$

$$\text{havoc}(s, \text{Vars}) \triangleq \forall v \in \text{Vars} \bullet s[v \mapsto z]$$

where  $z$  is a fresh variable (implicitly  $\exists$ -quantified).

<sup>7</sup> Note that the rule described in line 3 is slightly different from the one described in Sec. 3 because no consistency check is performed. Instead, the consistency check is postponed and done by the first base case at line 2.

`modifies( $\ell_1 \rightarrow \dots \rightarrow \ell_n$ )` takes a sequence of transitions and returns the set of variables that may be modified during its symbolic execution.

Intuitively, `invariant` clears the symbolic store of all variables modified in the loop (using the `havoc` function) and then enhances the path condition  $\Pi$  of the symbolic state with the invariants from the abstract interpreter.

Let us now explain the *grey* boxes in Fig. 3. Lines 5-7 in `BackwardDepsV` cover the case when a loop header has been encountered. The main purpose here is to abstract the current symbolic state by using the loop invariant obtained from the abstract interpreter. The algorithm calls the function `invariant` (at line 6) with the transitions in the loop so as to obtain a copy of the current symbolic state annotated with the approximate loop invariant in its path condition. At line 7, the `UnwindTreeV` procedure is called on the resulting symbolic state to explore the symbolic subtree of the loop.

If the symbolic execution encounters a loop backedge (lines 8-10) from  $\ell$  to  $\ell'$  it halts and backtracks. The reason is that the loop header at  $\ell'$  has already been symbolically executed with a loop invariant. Hence there is no need to continue the loop since the invariant ensures that no new feasible paths will be encountered if it is explored again. This is our basic mechanism to make the symbolic execution of the loop finite.

Finally, the main algorithm to handle loops, `BackwardDepsLoopV`, makes calls to the function `BackwardDepsV` until there is no change detected in the symbolic state of any program point. We present it in its simplest form, but it can be easily optimized to call `BackwardDepsV` only with the loop in which the change was detected.

## 5 Results

We implemented the path-sensitive slicer described in this paper and performed experiments to address the following questions:

1. Is our path-sensitive slicer practical for medium-size programs?
2. What is the impact of *reusing*?
3. How effective is a path-sensitiveness slicer against a path-insensitive version?

Our proof-of-concept implementation models the heap as an array. A flow-insensitive pointer analysis is used to partition updates and reads into alias classes where each class is modeled by a different array. Given an operation that involves pointers the sets *def* and *use* utilize the results of the pointer analysis. For instance, given the statement `*p = *q` the set *def* contains everything that might be pointed to by `p` and the set *use* includes everything that might be pointed by `q`. A theorem prover is used to decide linear arithmetic formulas over integer variables and array elements in order to check the satisfiability and entailment of formulas, and computing interpolants and witnesses. Program are first annotated with approximate loop invariants using the abstract interpreter `InterProc` [15]. Functions are inlined and external functions are modeled as having no side effects and returning an unknown value.

We used several instrumented device driver programs previously used as software model checking benchmarks: `cdaudio`, `diskperf`, `floppy`, and `serial`. In addition, we also considered `mpeg`, the `mpeg-1` algorithm for compressing video, and `fcron.2.9.5`, a cron daemon. For the slicing criterion we consider variables that may be of interest during

Program	LOC	Path-Insens		Path-Sens		
		Size Red	Time	Reuse		No Reuse
				Size Red	Time	Time
mpeg	5K	4%	21s	8%	628s	$\infty^1$
diskperf	6K	32%	2s	57%	94s	$\infty$
floppy	8K	36%	9s	47%	263s	$\infty$
cdaudio	9K	23%	10s	52%	301s	$\infty$
serial	12K	39%	16s	50%	395s	$\infty$
fcron.2.9.5	12K	42%	32s	61%	832s	$\infty$
Mean		<b>23%</b>	<b>15s</b>	<b>38%</b>	<b>418s</b>	–

**Table 1.** Results on Intel 3.2Gz 2Gb. <sup>1</sup> timeout after 2 hours or 2.5 Gb of memory consumption debugging tasks. For the instrumented software model checking programs, we choose as the slicing criterion the set of variables that appear in the safety conditions used for their verification in [10]. In the case of mpeg we choose a variable that contains the type of the video to be compressed. Finally, in fcron.2.9.5 we choose all the file descriptors opened and closed by the application.

Table 1 compares our path-sensitive slicer (columns labelled with Path-Sens) against the same slicer but without path-sensitivity (labelled with Path-Insens). Path-insensitivity is achieved by the following modifications in our slicer: (1) considering all paths as feasible, and (2) always forcing reuse. These changes have the same effect as always merging the abstract states along incoming edges in a control-flow merging node. In other words, they mimic running a path-insensitive slicer on the original CFG. We could have used a faster off-the-shelf path-insensitive program slicer (using e.g., [11]), however, our objective here is to isolate the impact of path-sensitivity and hence, we decided to perform the comparison on a common platform to produce the fairest results. Finally, we also tried running with different abstract domains, such as octagons and polyhedra, to generate loop invariants and the results were the same.

The column LOC represents the number of lines of program without comments. The column Size Red shows the reduction in slice size (in %) wrt the original program size. The reduction size is computed using the formula  $(1 - \frac{\text{size of slice}}{\text{size of original}}) \times 100$ . By *size* we mean all executable statements in the program, excluding type declarations, unused functions, comments, and blank lines. A minor complication here is that the SET may contain multiple *instances* of program points in the CFG, as can be seen in Fig. 1(c). To compare the reduction in slice sizes fairly, we use the rule mentioned at the beginning of Sec. 4 to compute slices: a transition in the original CFG is included in the slice if any of its instances in the SET is included in the slice.

The column Time reflects the running time of the analysis in seconds excluding the alias analysis and the external abstract interpreter. Column Reuse is our path-sensitive slicer with reusing, and No Reuse uses the same symbolic execution engine with automatic loop invariants but without interpolation and witness paths. Finally, we summarize in row Mean the numbers of columns Size Red and Time by computing their *geometric* and *arithmetic* mean, respectively.

We summarize our results as follows. The running times (column Reuse) of our path-sensitive slicer (with a mean of 418 secs) are reasonable considering the size of the programs and the current status of our prototype implementation which can be optimized significantly. The analysis of mpeg is especially slow and it is due to the existence

of many nested loops. On the other hand, the reuse of solutions clearly pays off. Without our reuse mechanism (column No Reuse) we were not able to finish the analysis of any program after a timeout of 2 hours or memory consumption of 2.5 Gb. Finally, the improvement in terms of reduction shown in column Reuse is roughly 38% against only 23% of its path-insensitive counterpart (column Path-Insens). Again, the mpeg program is an exception since the size of the slices in both Path-Insens and Path-Sens are quite big (i.e., very small reduction). The reason is that in mpeg all the computations depend on the type of video to be compressed which is our slicing criterion.

## 6 Related Work

Static slicing remains a very active area of research. We limit our discussion to the most relevant works that take into account path-sensitiveness. We also discuss pruning techniques that might have influenced our work.

**Fully path-sensitive methods.** Conditioned slicing described in [3, 5, 6] performs symbolic execution in order to exclude infeasible paths before applying a static slicing algorithm, so they are fully path-sensitive (for loop-free programs) similar to us. However, they perform full path enumeration and essentially explore the search space of the naive SET. Hence, they suffer from the path explosion problem.

**Partially path-sensitive methods.** A more scalable but not fully path-sensitive approach is described by Snelting et al. [20, 17, 19]. They compute the dependency between two program points  $y$  and  $x$  using the Program Dependence Graph (PDG) [11] and apply the following rule to remove spurious dependencies:  $I(y, x) \Rightarrow \exists \bar{v} : PC(y, x)$ , where  $I(y, x)$  stands for  $y$  *influences*  $x$  (i.e., there is a dependency at  $x$  on  $y$ ),  $\bar{v}$  is some assignment of values to program variables and  $PC(y, x)$  is the path condition from  $y$  to  $x$ . Essentially it means that if the path condition from  $y$  to  $x$  is found to be unsatisfiable, then there is definitely no influence from  $y$  to  $x$ . If there are multiple paths between two points, the path condition is computed as a disjunction of each path.

For the program in Fig. 1(a), [20, 17, 19] would proceed as follows. In the PDG there will be a dependency edge from  $\ell_8$  to  $\ell_1$ , hence they would check to see if the path condition  $PC(1, 8)$  is unsatisfiable. First they calculate the path condition from  $\ell_4$  to  $\ell_8$  as  $PC(4, 8) \equiv (x = 1 \wedge y > 0 \wedge z = x) \vee (y = 0 \wedge y > 0 \wedge z = x) \equiv (x = 1 \wedge y > 0 \wedge z = x)$ . Now they use this to calculate  $PC(1, 8) \equiv (x = 0 \wedge y = 5 \wedge ((a > 0 \wedge b = x + y \wedge PC(4, 8)) \vee (a \leq 0 \wedge PC(4, 8))))^8$ , which is not unsatisfiable. Hence the statement  $x:=0$  at  $\ell_1$  will be included in the slice.

The fundamental reason for this is that for [20, 17, 19], path conditions are only necessary and not sufficient, so false alarms in examples such as the above are possible. An important consequence of this is the fact that even for loop-free programs, their algorithm cannot be considered “exact” in the sense described in Sec. 1. However, our algorithm guarantees to produce no false alarms for such programs.

Finally, another slicer that takes into account path-sensitiveness up to some degree is Constrained slicing [8] which uses *graph rewriting* as the underlying technique. As the

<sup>8</sup> We have simplified this formula since [20, 17, 19] uses the SSA form of the program and adds constraints for  $\Phi$ -functions, but the essential idea is the same.

graph is rewritten, modified terms are tracked. As a result, terms in the final graph can be tracked back to terms in the original graph identifying the slice of the original graph that produced the particular term in the final graph. The rules described in [8] mainly perform constant propagation and dead code detection but not systematic detection of infeasible paths. More importantly, [8] does not define rules to prune the search space.

**Interpolation and SAT.** Interpolation has been used in software verification [1, 10, 16, 12] as a technique to eliminate facts which are irrelevant to the proof. Similarly, SAT can explain and record failures in order to perform conflict analysis. By traversing a reverse implication graph it can build a *nogood* or conflict clause which will avoid making the same wrong decision. Our algorithm has in common the use of some form of *nogood learning* in order to prune the search space. But this is where the similarity ends. A fundamental distinction is that in program verification there is no solution (e.g., backward dependencies) to compute and hence, there is no notion of reuse and the concept of witness paths does not exist. [9] uses interpolation-based model checking techniques to improve the precision of dataflow analysis but still for the purpose of proving a safety property.

Finally, a recent work of the authors [13] has been a clear inspiration for this paper. [13] uses interpolation and witnesses as well to solve not an analysis problem, but rather, a *combinatorial optimization* problem: the Resource-Constrained Shortest Path (RCSP) problem. Moreover, there are other significant differences. First, [13] is totally defined in a finite setting. Second, [13] considers only the narrower problem of extraction of bounds of variables for loop-free programs while we present here a general-purpose program analysis like slicing. Third, this paper presents an implementation and demonstrates its practicality on real programs.

## 7 Conclusions

We presented a path-sensitive backward slicer. The main result is a symbolic execution algorithm which excludes infeasible paths while pruning redundant paths. The key idea is to halt the symbolic execution while reusing dependencies from other paths if some conditions hold. The conditions are based on a notion of interpolation and witness paths at aiming to detect whether the exploration of a path can improve the accuracy of the dependencies computed so far by other paths. We have demonstrated the practicality of the approach with a set of real C programs. Finally, although this paper targets slicing our approach can be generalized and applied to other backward program analyses providing them path-sensitiveness.

## References

1. T. Ball, R. Majumdar, T. Millstein, and S. K. Rajamani. Automatic predicate abstraction of C programs. In *PLDI'01*, pages 203–213.
2. Leeann Bent, Darren C. Atkinson, and William G. Griswold. A comparative study of two whole program slicers for c. Technical report, University of California at San Diego, La Jolla, CA, USA, 2001.
3. Gerardo Canfora, Aniello Cimitile, and Andrea De Lucia. Conditioned program slicing. *Information and Software Technology*, 40, no. 11-12:595–607, 1998.
4. W. Craig. Three uses of Herbrand-Gentzen theorem in relating model theory and proof theory. *Journal of Symbolic Computation*, 22, 1955.
5. Sebastian Danicic, Chris Fox, and Chris Harman. Consit: A conditioned program slicer. In *ICSM'00*, pages 216–226.
6. M. Daoudi, L. Ouarbya, J. Howroyd, S. Danicic, M. Harman, C. Fox, and M.P. Ward. Consus: A scalable approach to conditioned slicing. *Working Conference on Reverse Engineering*, 2002.
7. E. W. Dijkstra. *A Discipline of Programming*. Prentice-Hall Series in Automatic Computation. Prentice-Hall, 1976.
8. John Field, G. Ramalingam, and Frank Tip. Parametric program slicing. In *POPL '95*, pages 379–392.
9. J. Fischer, R. Jhala, and R. Majumdar. Joining dataflow with predicates. In *ESEC/FSE-13*, pages 227–236, 2005.
10. T. A. Henzinger, R. Jhala, R. Majumdar, and K. L. McMillan. Abstractions from proofs. In *31st POPL*, pages 232–244. ACM Press, 2004.
11. S. Horwitz, T. Reps, and D. Binkley. Interprocedural slicing using dependence graphs. In *PLDI '88*, pages 35–46.
12. J. Jaffar, J.A. Navas, and A. Santosa. Unbounded Symbolic Execution for Program Verification. In *RV'11*, 2011.
13. J. Jaffar, A. E. Santosa, and R. Voicu. Efficient memoization for dynamic programming with ad-hoc constraints. In *23rd AAI*, pages 297–303. AAAI Press, 2008.
14. J. Jaffar, A. E. Santosa, and R. Voicu. An interpolation method for CLP traversal. In *15th CP*, volume 5732 of *LNCS*. Springer, 2009.
15. G. Lalire, M. Argoud, and B. Jeannet. The Interproc Analyzer. <http://pop-art.inrialpes.fr/people/bjeannet/bjeannet-forge/interproc>, 2009.
16. K. L. McMillan. Lazy annotation for program testing and verification. In *22nd CAV*, 2010.
17. Torsten Robschink and Gregor Snelting. Efficient path conditions in dependence graphs. In *ICSE '02*, pages 478–488.
18. S Seo, H Yang, and K Yi. Automatic construction of hoare proofs from abstract interpretation results. In *APLAS'03*, pages 230–245.
19. G Snelting, T Robschink, and J Krinke. Efficient path conditions in dependence graphs for software safety analysis. volume 15, pages 410–457.
20. Gregor Snelting and Abteilung Softwaretechnologie. Combining slicing and constraint solving for validation of measurement software. In *SAS*, pages 332–348, 1996.
21. M. Weiser. Program slicing. In *ICSE '81*, pages 439–449, 1981.