

Differential Assertion Checking

Shuvendu K. Lahiri
Microsoft Research
Redmond, WA, USA

Kenneth L. McMillan
Microsoft Research
Redmond, WA, USA

Rahul Sharma
Stanford University
CA, USA

Chris Hawblitzel
Microsoft Research
Redmond, WA, USA

ABSTRACT

Previous version of a program can be a powerful enabler for program analysis by defining new relative specifications and making the results of current program analysis more relevant. In this paper, we describe the approach of *differential assertion checking* (DAC) for comparing different versions of a program with respect to a set of assertions. DAC provides a natural way to write relative specifications over two programs. We introduce a novel modular approach to DAC by reducing it to safety checking of a composed program, which can be accomplished by standard program verifiers. In particular, we leverage automatic invariant generation to synthesize relative specifications for pairs of loops and procedures. We provide a preliminary evaluation of a prototype implementation within the SYMDIFF tool along two directions (a) soundly verifying bug fixes in the presence of loops and (b) providing a knob for suppressing alarms when checking a new version of a program.

Categories and Subject Descriptors

D.2.4 [SOFTWARE ENGINEERING]: Software/Program Verification—*Assertion checkers, Formal methods*

General Terms

Verification, Reliability

Keywords

differential analysis, verification, regressions

1. INTRODUCTION

There are several factors limiting the adoption of static analysis tools in the hands of developers. For static assertion checking, these include the need to define an assertion (or specification) to check, to provide environment specifications and to provide auxiliary invariants for loops and

procedures. Although many auxiliary invariants can be synthesized automatically by invariant generation methods, the undecidable nature (or the high practical complexity) of assertion checking precludes complete automation for a general class of user-supplied assertions.

It has often been proposed that utilizing previous versions of an evolving program can significantly reduce the cost of program analysis [24]. Such approaches run in two primary directions. First, in the presence of program refactoring, two versions can be checked for semantic equivalence to ensure the correctness of the transformation [25, 12, 19]. Second, verification can be performed *incrementally*, for example by carrying over invariants that are unaffected by the syntactic changes [28]. Although these techniques are useful in their own right, they are applicable in very limited contexts. First, most software changes (including some called refactoring) induce some behavioral change. Equivalence checking is too strong for such cases. Moreover, incremental verification can only be performed when the previous version does not have any false warnings — unfortunately, this is too strong a requirement for the usage of static analysis tools. Such false warnings have to be either removed by manually specifying additional invariants, or the tool has to resort to ad-hoc heuristics to suppress a class of warnings. The former seriously undermines the productivity gained from the use of static analysis, whereas the latter leads to brittle tools that may suppress true bugs.

In this paper, we propose another direction for exploiting previous versions of a program as an implicit specification, which appears to open up an interesting space for trading off soundness for cost required to apply an assertion checker. Our observation is simple:

We can often prove *relative correctness* between two similar programs with respect to a set of assertions statically with significantly lower cost than ensuring absolute correctness.

Given a program P with a set of assertions A , one traditionally asks whether there is an environment for P in which one of the assertions in A fails. One can instead ask a relative version of this question: given two versions P and P' containing a set of assertions A , does there exist an environment in which P passes but P' fails? We formalize this idea as the problem of *differential assertion checking* (DAC) — checking two versions of a program with respect to a set of assertions. Although this provides a weaker guarantee of correctness of P' , it closely corresponds to an interesting

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ESEC/FSE '13, August 18–26, 2013, Saint Petersburg, Russia
Copyright 2013 ACM 978-1-4503-2237-9/13/08 ...\$15.00.

class of bugs (*regressions*) that are often most relevant to a developer and have a good chance of getting fixed. Moreover, we argue that DAC has several desirable traits, when checking absolute correctness is rife with false alarms:

1. DAC allows for writing natural relative specifications without a lot of modeling (additional ghost variables) to express the properties.
2. DAC can be used to show that bug fixes do not cause additional regressions for a set of assertions.
3. Exploiting the structural similarity of programs P and P' allows simple relative specifications to answer the relative questions.

An idea similar to DAC was earlier proposed in the context of filtering false alarms for concurrent programs [17] (we discuss subtle differences in Section 3). At a high level, one can see this work as applying the idea towards evolving programs and extending the idea to deal with unbounded loops and recursion ([17] was restricted to bounded programs).

1.1 Motivating example

<pre>void StringCopy.1(wchar_t *dst, wchar_t *src, int size) { wchar_t *dtmp = dst, *stmp = src; int i; for (i = 0; *stmp && i < size - 1; i++) *dtmp++ = *stmp++; *dtmp = 0; }</pre>	<pre>void StringCopy.2(wchar_t *dst, wchar_t *src, int size) { wchar_t *dtmp = dst, *stmp = src; int i; for (i = 0; i < size - 1 && *stmp; i++) *dtmp++ = *stmp++; *dtmp = 0; }</pre>
--	--

Figure 1: Motivating example (in C): two versions of `StringCopy` (Figure 1 [15]).

Consider the two versions of the procedure `StringCopy` described in Figure 1. The version `StringCopy.2` is a procedure for copying the contents of a `char` buffer `src` into `dst`, described in an earlier work [15]. Let us first ignore `StringCopy.1`, which is a buggy version of `StringCopy.2`. In this paper, we adopt the convention that procedures on the left side of the figures correspond to buggy versions and those on the right correspond to correct versions. Let us illustrate the complexities of verifying the memory safety of `StringCopy.2` in isolation.

1. To specify memory safety, one needs to define the *bounds* of a buffer for C programs (unlike Java or C#). This can be accomplished by adding a ghost variable `Bound` that maps each allocated pointer (such as `dst`) to a non-negative integer. One possible way to specify the memory safety is to precede any dereference `*e` with the assertion `assert Bound(e) > 0`.
2. One needs a *precondition* that the bounds of `dst` and `src` have some relationship with `size`, and the two buffers are disjoint.

3. Finally, one needs to write a loop invariant to record that `dtmp` always points inside the buffer pointed by `dst`, among other things.

Even for such a simple procedure, specifying and verifying the memory safety can be quite complex if the user is left to define the assertions, environment conditions and intermediate invariants.

Now we define *relative* memory safety of `StringCopy.2` with respect to `StringCopy.1`. First, observe that the difference in the two versions lies in the loop exit condition where the conjunction (`&&`) is applied in reverse order — this gives different behaviors due to the short-circuit semantics of `&&`. We want to check that `StringCopy.2` accesses only the memory locations which `StringCopy.1` accesses for any input. We can define and check relative memory safety in a generic fashion as follows:

1. Define an uninterpreted predicate *Valid* that maps each pointer to a Boolean value. Each dereference `*e` is preceded by `assert Valid(e)`.
2. Let `ok.i` be a global Boolean variable for `StringCopy.i` procedure that is true if no assertion has failed. We replace `assert ϕ` by code that sets `ok.i` to false if `ϕ` is false. We say `StringCopy.2` is correct relative to `StringCopy.1` if, when both start in the same state (parameters and the heap) and both terminate, if the former terminates in a state where `ok.1` is true, then the latter also terminates in a state satisfying `ok.2`.
3. Assuming the two loops are automatically extracted as tail-recursive procedures (§ A) `loop.1` and `loop.2` respectively, we show how to construct a composed procedure for the two loops and attach a simple relative specification on the composed procedure.

```
pre  stmp.1 == stmp.2 &&
     dtmp.1 == dtmp.2 &&
     Mem_char.1 == Mem_char.2 &&
     i.1 == i.2 &&
     size.1 == size.2 &&
     ok.1 <==> ok.2

post ok.1 ==> ok.2 &&
     dtmp.1 == dtmp.2
proc MS_loop.1_loop.2(dst.1, ..., dst.2, ...);
```

Here `pre` refers to a precondition and `post` refers to a postcondition, and `Mem_char.i` refers a global array that models the state of the heap. Moreover, we show how such a specification can be inferred using the techniques in this paper.

Note that we did not require any precondition about the inputs to the program, nor any correlation about the bounds nor any relationship with null-terminated buffers. This checking succeeds and we have proven that `StringCopy.2` has a memory footprint no larger than `StringCopy.1`. On the other hand, if one were to check the relative correctness of `StringCopy.1` with respect to `StringCopy.2` under the relative memory safety specification, one would get a counterexample where `size` equals 0 and pointer `src` does not satisfy `Valid`. This counter-example captures the seeded bug: an address that `StringCopy.1` dereferences but `StringCopy.2` does not.

x	\in	$Vars$	
R	\in	$Relations$	
U	\in	$Functions$	
e	\in	$Expr$	$::= x \mid c \mid U(e, \dots, e) \mid \text{old}(e)$
ϕ	\in	$Formula$	$::= \text{true} \mid \text{false} \mid e \text{ relop } e \mid \phi \wedge \phi \mid \neg\phi \mid R(e, \dots, e) \mid \dots$
s	\in	$Stmt$	$::= \text{skip} \mid \text{assert } \phi \mid \text{assume } \phi \mid x := e \mid \text{havoc } x \mid s; s \mid x := \text{call } f(e, \dots, e)$
c	\in	$CFStmt$	$::= L : \mid \text{goto } L_1, \dots, L_n \mid \text{return}$
f	\in	$Body$	$::= c \mid s; f \mid f; f$
p	\in	$Proc$	$::= \text{int } f(x_f : \text{int}, \dots) : r_f \{ f_{body} \}$

Figure 2: A simple programming language. The set of goto statements do not form any cycles in the control flow graph.

1.2 Overview

In the rest of the paper, using the background developed in § 2, we formalize the notion of differential assertion checking (DAC) (§ 3), and illustrate its use for defining relative specifications (§ 3.1). We provide an algorithm for checking DAC modularly by transforming the relative correctness problem into verifying assertions over a single composed program (§ 4). This allows us to leverage standard off-the-shelf program verifiers and invariant generation tools to check the relative correctness problem. We demonstrate a simple scheme based on HOUDINI [11] that suffices for a class of programs (§ 5.1). We have created a prototype implementation of our method inside SYMDIFF [19], a semantic differencing tool. We evaluate the tool along two different directions. First, we use DAC to soundly verify that the version after a bug fix is relatively correct with respect to the buggy version (§ 6.1). Second, we show that DAC can provide a systematic knob for suppressing alarms when analyzing a new version of a program (§ 6.2). Together, the experiments indicate the potential of DAC to be a generic framework to exploit previous versions of a program.

2. BACKGROUND

Figure 2 describes a simple programming language (a subset of the BOOGIE [2] programming language) with recursive procedures and an assertion language. We assume that loops are already desugared into recursive functions of this language (we describe a method in § A). The language supports variables ($Vars$) and various operations on them. Expressions ($Expr$) can be variables, constants, or the result of applying a (possibly interpreted) function U to a list of expressions. The expression $\text{old}(e)$ refers to the value of e at the entry to a procedure. $Formula$ represents Boolean valued expressions and can be the result of (interpreted or uninterpreted) relational operations on $Expr$, Boolean operations ($\{\wedge, \neg\}$), or possibly quantified expressions ($\forall u : \text{int}.\phi$). Note that the programming language is fairly expressive and can be used to model arrays. An array can be modeled in this language, by introducing two special functions $sel \in Functions$ and $upd \in Functions$; $sel(e_1, e_2)$ selects the value of a map value e_1 at index e_2 , and $upd(e_1, e_2, e_3)$ returns a new map value by updating a map value e_1 at location e_2 with value e_3 .

A state of a program at a given program location is a valuation of the variables in scope (procedure parameters,

locals and global variables) and a program counter pc that indicates the next statement to be executed. A program consists of a set of *basic blocks*, where each basic block consists of a statement $s \in Stmt$ terminated with a control flow statement $CFStmt$ (*goto* or *return* statement). A *goto* statement $\text{goto } L_1, \dots, L_n$ non-deterministically sets the pc to any one of the n labels. We restrict the use of *goto* statements to not form any cycles in the control flow graph. The statement *skip* denotes a no-op. The statement *assert* ϕ is used to statically check that the formula ϕ holds; *assert* ϕ has no effect on the dynamic state. The statement *assume* ϕ behaves as a *skip* when the formula ϕ evaluates to *true* in the current state; else the execution of the program is *blocked*. The assignment statement is standard, *havoc* x scrambles the value of a variable x to an arbitrary value, and $s; t$ denotes the sequential composition of two statements s and t . Conditional statements are modeled by using the *goto* statement and *assume* statements. Procedure calls are denoted using the *call* statement, and can have a side effect by modifying one of the global variables.

Let Σ be the set of all states for a program. For any procedure $p \in Proc$, we assume a transition relation $\mathcal{T}_p \subseteq \Sigma \times \Sigma$ that characterizes the input-output relation of the procedure p . In other words, two states $(\sigma, \sigma') \in \mathcal{T}_p$ if there is an execution of the procedure p starting at σ and ending in σ' . The transition relations can be defined inductively on the structure of the program and is fairly standard for our simple language [2].

There are a host of tools for modeling most high level languages (such as C, C#, Java) in this language (such HAVOC [7] for C). We note that such translations use arrays to model the heap (e.g. an array per field in Java) where the arrays are indexed by objects or pointers. We defer further discussion of the translations to these earlier works.

3. DAC

In this section, we formalize our approach of differential assertion checking (DAC). The basic concept of DAC appears in a previous work in the context of filtering false alarms in verification of concurrent programs using sequential executions [17]. However, it was described in a simpler setting of bounded programs: loops were unrolled and recursive procedures were inlined a bounded number of times.

Before proceeding, we establish a few notations that we follow in the paper unless explicitly stated otherwise. First, we assume that any assertion *assert* ϕ is replaced by the assignment $\text{ok} := \text{ok} \wedge \phi$ to a global ok variable. Second, given that we are considering two versions P_1 and P_2 of a program, we suffix the names of procedures, globals (including ok) and parameters with the version number. Third, we label a state σ as *failing* if ok variable is *false* in σ . Finally, we assume a one-one (not necessarily onto) mapping between the globals, procedures, and their parameters between the two versions; we often equate states from two versions when we really mean that the two states assign the same value to the mapped variables of the two states.

Definition 1. (Differential assertion checking) Given two procedures p_1 and p_2 , p_2 has a differential error with respect to p_1 (denoted as $DAC(p_2, p_1)$) if there exists an input state σ such that (1) there exists a state σ'_1 such that $(\sigma, \sigma'_1) \in \mathcal{T}_{p_1}$ and σ'_1 is non-failing, and (2) there exists a state σ'_2 such that $(\sigma, \sigma'_2) \in \mathcal{T}_{p_2}$ and σ'_2 is failing.

We define a procedure p_2 to be *relatively correct with respect to* p_1 if $DAC(p_2, p_1)$ does not hold.

The above definition differs from the definition of *differential error* ($DiffErr(p_2, p_1)$) [17] in a subtle way. The difference lies in whether we insist the input σ to be non-failing for *every* execution in p_1 (in $DiffErr(p_2, p_1)$) as opposed to being failing on some execution in p_1 (in $DAC(p_2, p_1)$). We provide a simple example that distinguishes the two views in Figure 3.

<pre>proc p1() { havoc x; if (x) assert false; }</pre>	<pre>proc p2() { assert false; }</pre>
--	--

Figure 3: Example differentiating $DiffErr$ and DAC

For this example, $DAC(p_2, p_1)$ holds as there is a state (empty) from which p_1 succeeds (when the internal variable x is assigned `false`) and p_2 fails. However, $DiffErr(p_2, p_1)$ does not hold because there is no input state from which all executions are non-failing for p_1 . It is easy to observe that if $DiffErr(p_2, p_1)$ holds then $DAC(p_2, p_1)$ holds, but not otherwise.

The definition of $DiffErr$ was motivated by comparing concurrent interleaved executions with their sequential counterparts. We adopt the slightly modified definition for DAC to several reasons. First, the check for $DAC(p_2, p_1)$ can be encoded very naturally using single program verifiers:

<pre>assume i1 == i2 && g1 == g2; call p1(i1); call p2(i2); assert (ok.1 ==> ok.2);</pre>
--

where we use i and g to denote parameters and globals. On the other hand, the $DiffErr$ check is more complicated because checking it is undecidable even for bounded programs. This added complexity is not needed for comparing similar versions of a program; we have found that internal non-deterministic choices are less common. Whenever non-determinism is present (say reading chars using `scanf`), the choices can be aligned on the two sides to return the same arbitrary sequence of choices in the two programs (see [19]). In such a modeling, the non-deterministic choices become reads from an input array, thereby converting internal non-determinism to input non-determinism.

3.1 Relative specifications

Recall that writing meaningful specifications often require access to a host of ghost state that is not present explicitly as part of the program state (§ 1.1). In addition to checking existing assertions in the two versions differentially, DAC also facilitates writing relative specifications using the same syntax of single program assertions. Instead of defining the buffer overrun checks on the two programs and checking them differentially, it often helps to pose questions such as: are there inputs for which P_2 accesses buffer regions that are not accessed by P_1 ? Such specifications can be written by introducing an uninterpreted predicate $Valid$ and adding an assertion before accessing any pointer p : `assert Valid(p)`. Such a specification will be useless for checking a single program (every pointer dereference might be flagged as a warning), but will naturally provide a relative specification.

Moreover, such a specification can be strengthened using semantics of the particular property that is desired. For example, when checking for non-null pointer dereferences, one can constrain the predicate by adding an axiom:

$$\text{axiom}(\forall x : \text{int} :: x \neq 0 \Rightarrow \text{Valid}(p))$$

Similarly, while checking for buffer overflows, one can add an axiom:

$$\text{axiom}(\forall x : \text{int}, y : \text{int} :: x \leq y \Rightarrow \text{Valid}(y) \Rightarrow \text{Valid}(x))$$

This will allow the DAC to not show a warning when the program P_2 accesses an index that is smaller than an index accessed by P_1 . This is specially useful when the entire history of indices accessed by P_1 is not stored (especially while doing a modular proof of DAC (§ 4) that only records an abstraction of the history of accesses on the two programs). Finally, one can even capture properties such as equivalence of two procedures (modulo termination). For a procedure $p \in P$, let o be the set of out parameters and g be the set of globals modified by p . If we assert $ValidEQ(o, g)$ (for an uninterpreted predicate $ValidEQ$) on the post-state of p and then perform DAC on two versions p_1 and p_2 , then the relative specification is correct if and only if the two programs are equivalent.

4. MODULAR DAC

In the previous section, we defined the problem of differential assertion checking $DAC(p_2, p_1)$ for a pair of procedures p_1 and p_2 . In this section, we provide a mechanism to check for $DAC(p_2, p_1)$ (or rather verify that p_2 is relatively correct with respect to p_1) in a procedure modular manner. In other words, we will verify the relative correctness without inlining the callers inside a procedure, but rather using some specifications. We provide a program transformation technique that compiles the relative correctness check of two programs P_1 and P_2 into a single composed program, which can be analyzed by an off-the-shelf program verifier. In particular, the transformation allows us to leverage existing invariant inference mechanisms for single programs for inferring relative specifications. The transformation is not specific to the problem of differential assertion checking, and is applicable whenever there is a need to compare two programs.

4.1 Composed program

<pre>proc f1(x1): r1 modifies g1 { s1; L1: w1 := call h1(e1); t1 }</pre>	<pre>proc f2(x2): r2 modifies g2 { s2; L2: w2 := call h2(e2); t2 }</pre>
--	--

Given two programs P_1 and P_2 each containing a set of procedures, and a one-one mapping between procedures, let us consider two particular mapped procedures $f1 \in P_1$ and $f2 \in P_2$. We have specified the modified set of globals for each procedure using `modifies` keyword. For ease of exposition, we have assumed that the read set of a procedure is a superset of the set of modified variables.

```

proc MS_f1_f2(x1,x2) returns (r1,r2)
modifies g1, g2
{
  // initialize call witness variables
  b_l1, b_l2, ... := false, false, ...;

  [[s1 :]]
L1:
  i_l1, gi_l1 := e1, g1 ; //store inputs
  call w1 := h1(e1);
  b_l1 := true; //set call witness
  o_l1, go_l1 := w1, g1; //store outputs

  [[t1 :]]

  [[s2 :]]
L2:
  i_l2, gi_l2 := e2, g2 ; //store inputs
  call w2 := h2(e2);
  b_l2 := true; //set call witness
  o_l2, go_l2 := w2, g2; //store outputs

  [[t2 :]]

  //one block for each pair of call sites
  //for a pair of mapped procedures
  ....
  if (b_l1 && b_l2) { //for (L1,L2) pair
    //store the globals
    st_g1, st_g2 := g1, g2;

    g1, g2 := gi_l1, gi_l2 ;
    call k1, k2 := MS.h1.h2(i_l1, i_l2 );
    assume (k1 == o_l1 && g1 == go_l1);
    assume (k2 == o_l2 && g2 == go_l2);

    //restore globals
    g1, g2 := st_g1, st_g2;
  }
  ...
return;
}

```

Figure 4: Composed procedure for f1 and f2.

Figure 4 describes a composed procedure `MS_f1_f2` that is constructed for each pair of mapped procedures. First, note that the signature (parameters, modifies sets) of the procedure is a disjoint union of the signatures of the individual procedures. The body of `MS_f1_f2` consists of sequential composition of the bodies of `f1` and `f2`, in addition to some extra instrumentation. Since loops are already extracted as tail-recursive procedures, the body of any procedure contains no loops.

The instrumentations consist of two parts. The first part consists of storing the input and the output state at each call site. The second part consists of constraining the outputs of pairs of call sites (from different programs) to be the result of executing the corresponding composed procedure over the input states at the two call sites. This allows us to infer facts about pairs of procedure calls and to apply them in context.

We describe each of the steps in detail with respect to a pair of call sites from `f1` and `f2` respectively. At a given call site (say for label L1), we store the arguments and the

input value of global variables into local variables (`i_l1` and `gi_l1`) respectively. Since `f1` only modifies globals from `g1`, it suffices to store this subset of globals. Similarly, we record the returned value and the globals after return into local variables (`o_l1` and `go_l1`) respectively. Each call site also has a local Boolean *witness* variable (`b_l1`) that is initialized to `false` and set to `true` after the call has returned. The figure shows the transformation of the two particular call sites; other call sites in the remainder of the procedures are similarly instrumented (indicated by the double brackets in “[[si:]]).

After the instrumentation of the bodies of the two procedures, we add a conditional block for each pair of mapped call sites. The blocks are guarded by the Boolean witness variables for the call sites; these blocks are executed only when the corresponding call sites were encountered in an execution and both returned. Each block first stores the values of the globals into local `st_gi` variables. Next, it calls the composed procedure `MS_h1_h2` (this time for the pair of callees), with the calling contexts restored from the `gi_li` variables, passing stored arguments `i_li` as inputs to the composed procedure. The return values (returns and globals) are constrained to be the recorded values from after the two calls, using the `assume` statements. Finally, the globals are restored from the `st_gi` variables, erasing the effect of the call.

We use the notation $\sigma_1 \oplus \sigma_2$ to denote a composed state consisting of a state from the two programs with disjoint signatures.

THEOREM 1. *For two programs P_1 and P_2 and two procedure $p_1 \in P_1$ and $p_2 \in P_2$, $(\sigma_1, \sigma'_1) \in \mathcal{T}_{P_1}$ and $(\sigma_2, \sigma'_2) \in \mathcal{T}_{P_2}$ if and only if $(\sigma_1 \oplus \sigma_2, \sigma'_1 \oplus \sigma'_2) \in \mathcal{T}_{MS_{p_1-p_2}}$.*

PROOF. We only sketch the main ideas here. The first part of `MS_{p1-p2}` has the effect of executing `p1` and `p2` in parallel, recording the pre- and post-states of the procedure calls in ghost variables. The second part always has a terminating execution and has no effect. That is, by induction on recursion depth, we can assume the theorem for the call to `MS_{h1-h2}`. This guarantees a behavior for which the subsequent `assume` statements are true. Moreover the program’s global state is restored. Thus the net effect of `MS_{p1-p2}` is simply to execute `p1` and `p2`. \square

Theorem 1 illustrates that the transformation performed is not just limited to performing differential assertion checking, but provides a general method to exploit similarity between procedures in program proving. The main power of the transformation comes from providing the additional composed procedures over which one can write specifications towards the proof of a final specification (like DAC). An invariant inference engine now has the extra flexibility to infer invariants about the composed procedures in addition to the procedures in P_1 and P_2 .

Consider two versions of `Foo` (Figure 5) where the second version accesses fewer indices in the array `a`. Let us assume that the loops are extracted into procedures `Loop.1` and `Loop.2` respectively (omitted for brevity). Our approach will generate the following composed method `MS_Loop.1_Loop.2`. The relative specification (using the keyword `post`) says that if the values of `i` and `ok` are equal at the start of a loop execution, then `Loop.2` fails less often than `Loop.1`. This is an inductive specification, and also sufficient to prove the DAC property for the outer procedures `Foo.1` and `Foo.2`.

```

var a .1:[ int ]int;
const MAX: int;
proc Foo.1() {
  var i: int, t:int;
  i := 0;
  while (i <= MAX) {
    assert Valid(i);
    t := a.1[i];
    i := i + 1;
  }
}

```

```

var a .2:[ int ]int;
const MAX: int;
proc Foo.2() {
  var i: int, t:int;
  i := 0;
  while (i < MAX) {
    assert Valid(i);
    t := a.2[i];
    i := i + 1;
  }
}

```

Figure 5: Running example.

```

post (i.1 == i.2 && old(ok.1) <==> old(ok.2))
    ==> (ok.1 ==> ok.2)
proc MS_Loop.1.Loop.2(i.1, t.1, i.2, t.2)
  returns (i.1', t.1', i.2', t.2');
modifies ok.1, ok.2

```

The example also illustrates one other important aspect. The specifications of composed procedures typically have a simple relative form, but are not entirely trivial to obtain. If we had included the equality $t.1 == t.2$ alongside $i.1 == i.2$ our specification would have been too weak, since t is not initialized on entry to the loops. Mutual specifications are often mostly independent of the actual effect of procedures (a great advantage) but may not be the trivial equality over all the state variables in scope.

4.1.1 Relative vs. absolute specifications

On the other hand, let us consider the complexity of the specifications *without* the composed procedure. To prove the DAC property on the two versions of `Foo`, one will need to provide the following precondition for `Foo.2`:

```
pre forall j :: 0 <= j && j <= MAX ==> Valid(j)
```

Informally, this provides the weakest precondition of `Foo.1` to ensure that the procedure does not fail. In addition, we will need a loop invariant on `Loop.2` procedure:

```
pre 0 <= i.2 && i.2 <= MAX
```

Although this is another way to prove the DAC property, it demonstrates that one may require program specific (possibly quantified) invariants (since it talks about `MAX`) that may become arbitrarily complex to specify and more difficult to infer. On the other hand, the relative specification used for proving the DAC property using the composed procedure can be fairly easy to guess as it may depend little on details of the actual procedures.

5. INFERRING RELATIVE CONTRACTS

Since the composed procedures have the same syntax as the underlying procedures in P_i programs, we can use any invariant inference technique that can be used to generate invariants for P_i programs. In particular, we can use ideas based on abstraction interpretation [8], predicate abstraction techniques [13], and interpolants [22]. However, any invariant synthesis technique is necessarily incomplete and might either be limited by the underlying domain or may

diverge trying to find the inductive invariant. Therefore, it is wise to inject some domain knowledge while looking for invariants for proving differential properties like *DAC*.

In general, there are two forms of contracts for a composed procedure such as $MS_f_1_f_2$ in Figure 4. The precondition of such a procedure would be a predicate over the parameters and globals (i_1, i_2, g_1, g_2) , and the postcondition would be predicate over the input and output parameters and globals $(i_1, i_2, old(g_1), old(g_2), r_1, r_2, g_1, g_2)$ — we assume that the read sets are also included in g_i globals. Further, many natural two-state postconditions have the form $\phi(i_1, i_2, old(g_1), old(g_2)) \Rightarrow \psi(r_1, r_2, g_1, g_2)$. Finally, each of ϕ and ψ usually relate mapped variables (whenever such a mapping can be easily obtained by matching names or types) from the two programs using relations such as equality, inequality and Boolean implications.

5.1 Conjunctive relative specifications

We describe a simple scheme for synthesizing a subset of the specifications described above, namely conjunctive relative specifications. For each composed procedure we automatically generate a set of *candidate preconditions* and *candidate postconditions* and use the HOUDINI algorithm [11] to infer a subset of these that are inductive for the program and prove the specification. HOUDINI performs a greatest fix-point computation starting with the set of all candidate contracts as live (preconditions and postconditions) and kills a candidate when it cannot be proved modularly assuming the other live candidates. The process is repeated until either no candidate can be removed, or the desired specification can no longer be proved. In the former case, a sufficient inductive invariant has been synthesized for the specification; the latter case indicates either the property does not hold or the set of candidates is insufficient. For BOOGIE programs, one can use an efficient implementation of HOUDINI algorithm using the `/contractInfer` switch in BOOGIE [30].

Now, we describe the set of candidates that are automatically generated for each composed procedure such as $MS_f_1_f_2$ in Figure 4. For simplicity, we also assume that each program P_i has a single entry procedure (say p_i^0) that is not called from within P_i and all procedures in P_i have a body. For each f_i ($i \in \{1, 2\}$), let us denote I_i as in-parameters, M_i as the *ref* set of globals, R_i as the out-parameters and G_i as the mod set of globals. For each procedure other than the entry procedure, we first define the sets V_i as $I_i \cup M_i$ (for preconditions) and $R_i \cup G_i$ (for postconditions). For any pair of mapped variables $v_1 \in V_1$ and $v_2 \in V_2$, we add the following expressions as either preconditions or postconditions: (i) $\{v_1 \Rightarrow v_2, v_2 \Rightarrow v_1\}$ for Booleans, (ii) $\{v_1 \leq v_2, v_2 \leq v_1\}$ for integers and (iii) $\{v_1 = v_2\}$ otherwise. Given these candidates, HOUDINI algorithm generates the strongest inductive conjunctive invariant (if any) over these candidates that can prove the DAC specification.

6. EVALUATION

In this section, we describe an implementation and evaluation of DAC inside SYMDIFF [19]. SYMDIFF is an infrastructure for leveraging program verification techniques for comparing programs. The tool is agnostic to source languages (C, Java, C#, x86) as it operates on the BOOGIE intermediate verification language. It currently has a front-end for C programs (using the HAVOC [7] tool) that we use for our experiments. Internally, SYMDIFF leverages the efficient ver-

ification condition generation in BOOGIE [3] along with the Z3 [9] theorem prover to verify loop-free and call-free fragments. The implementation of DAC consists of around 800 lines of C# code and mainly performs the following program transformations: (i) introduces an `ok` variable and rewrites the assertions present in a program, (ii) generates the composed procedures (Figure 4), (iii) adds the DAC specification for the entry procedures, and (iv) generates the candidate contracts for the composed procedures (§ 5.1). In addition, for each procedure p , it adds a postcondition $ok \Rightarrow old(ok)$ — this captures the semantics that the `ok` variable can only transition from `true` to `false`.

In the next two subsections, we describe our experience with applying DAC towards two directions. First, we evaluate the inference of relative specifications for verifying bug fixes for a set of small C examples with unbounded loops (§ 6.1). Next, we evaluate the effectiveness of DAC as a mechanism for filtering alarms for evolving programs compared to checking assertions on a single program (§ 6.2).

6.1 Verifying bug fixes

Table 1 describes the result of performing DAC on a set of C examples (except `iter` which is a hand written BOOGIE example). Each example contains between one and three procedures with at least one loop. The first two examples are already described in this paper, `iter` in Section 4 and `strcpy` in Figure 1. The rest of the examples are drawn from the VERISEC suite containing “snippets of open source programs which contained buffer overflow vulnerabilities, as well as the corresponding patched versions.” [31]. For each of these benchmarks, we add an assertion `assert Valid(p)` before any dereference to a pointer expression p . This includes array accesses where $a[i]$ is treated as $*(a + n * i)$ for an array whose base type occupies n bytes. Performing DAC checks that the corrected version is dereferencing only the memory locations which the buggy version does and the bug fix has not inadvertently increased the memory footprint.

Example	# Glbs	# Cands	# Infrd
<code>iter</code>	2	13	6
<code>strcpy</code>	19	29	28
<code>apache-1</code>	23	88	72
<code>madwifi-1</code>	36	187	59
<code>madwifi-2</code>	30	141	117
<code>sendmail-1</code>	20	77	49
<code>sendmail-2</code>	24	65	56

Table 1: Bug fix verification results. “Glbs” denotes globals in the BOOGIE translation of each program, “Cands” denotes candidate preconditions or postconditions, “Infrd” denotes the subset of “Cands” that were inferred by HOUDINI.

The examples in the VERISEC suite range from around 20 to 50 lines of C code (see Figure 6 for the `sendmail-1` example). Table 1 indicates that the number of global variables is non-trivial in each example (except `iter` which is a manually encoded BOOGIE program). These globals (generated by HAVOC [7]) model various aspects of C semantics including maps for each pointer types and fields, allocation status of pointers, and deterministic sequence of values returned by functions such as `nondet_int` (Figure 6). Finding

the right relative specifications can be extremely time consuming given the sizes of product programs. Therefore, the inference is quite invaluable in discovering the relative invariants needed to prove the DAC property, even for these small C examples. Only one example (`apache-1`) required an additional (absolute) specification not generated by our tool — it specifies that a loop index variable never decreases. For rest of the benchmarks, we were able to automatically infer contracts which were sufficient to prove that the memory footprint of the correct version was no larger than the footprint of the buggy version.

The pair of procedures for `sendmail-1` in Figure 6 illustrates a couple of challenges for differential reasoning. First, note that the fix resets the counter `fb` to 0 under some condition. Therefore, the values of `fb` on the two programs will get out of sync after `fb` reaches `MAXLINE`, since the buggy program will continue to increment `fb`. Hence the precondition of the composed procedure for the loops only satisfies the specification $fb.2 \leq fb.1$ as opposed to $fb.2 == fb.1$. Second, if `Valid` is completely unconstrained, one may not be able to prove the DAC property modularly without using quantifiers in the invariants to record the history of accesses in the first loop. Instead, we constrain `Valid` by the axiom $\forall x, y :: x \leq y \wedge Valid(y) \Rightarrow Valid(x)$ (§ 3.1), allowing the simple relative specifications to prove the DAC property.

```

int main (void)
{
  ...
  fb = 0;
  while ((c1 = nondet_int ())
         != EOF) {
    /* OK */
    fbuf[fb] = c1;
    fb++;
    if (fb >= MAXLINE)
      fb = 0;
  }
  /* force out partial
   last line */
  if (fb > 0) {
    /* OK */
    fbuf[fb] = EOS;
  }
  return 0;
}

int main (void)
{
  ...
  fb = 0;
  while ((c1 = nondet_int ())
         != EOF) {
    /* BAD */
    fbuf[fb] = c1;
    fb++;
  }
  /* force out partial
   last line */
  if (fb > 0) {
    /* BAD */
    fbuf[fb] = EOS;
  }
  return 0;
}

```

Figure 6: Example of modular bug fix verification (`sendmail-1`). The “BAD” and “OK” denote buggy and fixed buffer accesses respectively.

Hence, we have demonstrated that DAC can be used for verification of bug fixes. Starting from buggy and correct versions of programs from a standard buffer overflow benchmark, DAC automatically infers relative contracts and proves that the bug fix does not introduce dereferences of new locations; hence, eliminating the possibility of a regression.

6.2 Filtering warnings

In this section, we evaluate the trade offs of differential reasoning as a mechanism for filtering warnings from a program verifier for evolving programs. When a single program is analyzed for some specification (say memory safety) by a verifier, for some programs, invariably there is a flood of

warnings. Many such warnings are false alarms due to the limitations of static checking. A developer in such a situation will need some knobs which can lead him to warnings of interest. In evolving software projects, a user is often less concerned with warnings that were present in the earlier releases.

In this section, we perform two case studies for exploring such knobs: with benchmarks from Software-artifact Infrastructure Repository [27, 10] and Windows device drivers¹. For this section, we check the DAC property with respect to the absence of null dereference errors. Each dereference of a pointer p is preceded with an assertion about $Valid(p)$. Unlike the previous section, we however do not solely focus on changes that correspond to bug fixes for this class of assertions.

For the purpose of this section, we have done several restrictions and simplifications. First, the loops present in any procedure is unrolled two times. This is done to separate the benefits of DAC from the precision gain obtained by using an invariant inference engine. Second, we only consider one candidate postcondition for the composed procedure where the mapped procedures are semantically equal. This is the default summary considered by SYMDIFF for performing equivalence checking. In other words, the summary of the composed procedure $MS_{p_1-p_2}$ is limited to either the procedure equivalence or the trivial summary **true**.

We instantiate the framework with five configurations: (i) **single**: each procedure in P_2 (without taking P_1 into account) affected by the change is checked modularly without any preconditions and callee postconditions. This is the default behavior of the static analysis performed by HAVOC. (ii) **sound**: when analyzing P_1 and P_2 differentially, we use the candidate summaries described above for the callees. (iii) **unsound**: we assume that callees do not modify the `ok` variables. This amounts to unsoundly assuming that callees do not fail even when called from different states in P_2 compared to P_1 . (iv) **shallow**: we unsoundly assume that callees are equivalent including the effect on the `ok` variables. (v) **nonmodular**: we check DAC non-modularly by inlining callees and do not use any specifications. We have designed the different options to compare modular DAC (represented closest by **sound**) with (a) non-differential reasoning (**single**), (b) non-modular DAC (**nonmodular**), (c) effect of increasing unsoundness (**unsound** and **shallow**), which in turn restricts the adversarial environments a static analysis can consider while analyzing a procedure, on a large class of examples. Note that the degree of unsoundness increases in going from **sound** to **unsound** to **shallow**. These modular analyses try to find a single input for an internal procedure for which P_1 does not fail, but P_2 does. On the other hand, **nonmodular** performs an analysis assuming equal inputs only for the entry procedures only and not for the internal procedures — hence it is incomparable with the other options. As expected, our experiments demonstrate that **sound** \supseteq **unsound** \supseteq **shallow** in terms of versions that have warnings. We have also observed that the runtime of **nonmodular** is often 10-100 times more expensive compared to the modular approaches.

Table 2 describes the results on the SIEMENS and SPACE suite of C benchmarks, available from the Software-artifact Infrastructure Repository [27]. Each program in this suite

¹Microsoft Windows Driver Kit (WDK), available at <http://www.microsoft.com/whdc/devtools/ddk/default.mspx>

has several versions (the column **versions**) that correspond to injecting various bugs encountered during the development of these benchmarks. However, these bugs are usually functional bugs (changing some conditional or mutating an arithmetic operation) that often do not manifest in null dereference errors. As can be seen from the table, the number of warnings (110 versions out of a total of 127 versions in 848 procedures) arising while checking null dereference absolutely (**single**) can be quite high, even when focusing on the procedures impacted by the change. In comparison, number of warnings progressively decreases with the use of **sound**, **unsound**, and **shallow** options. The **nonmodular** represents the true set of DAC errors; however, inlining does not scale to large programs such as SPACE. Out of these warnings, we have confirmed that 2 warnings in **schedule** are true null-dereference bugs caused by the change. We also notice that **unsound** and **shallow** options are very similar in nature, except that more procedures can fail in **unsound** (e.g. **schedule2**).

Figure 7 (from **schedule2**) shows an example where performing differential reasoning allowed suppressing a warning generated by **single**. On analyzing just a single procedure, every dereference of `job` is flagged as a warning, since the input value of `job` can be `null`. DAC is able to show relative correctness: the second program does not dereference a null pointer if the first one does not.

<pre> int get_process(prio, ratio, job) int prio; float ratio; struct process ** job; { ... <i>if(ratio < 0.0 </i> <i>ratio > 1.0)</i> <i>return(BADRATIO);</i> ... *job = *next; if(*job) { ... return(TRUE); } else return(FALSE); } </pre>	<pre> int get_process(prio, ratio, job) int prio; float ratio; struct process ** job; { ... <i>if(ratio < 0.0 </i> <i>ratio >= 1.0)</i> <i>return(BADRATIO);</i> ... *job = *next; if(*job) { ... return(TRUE); } else return(FALSE); } </pre>
--	---

Figure 7: Difference between “Single” and “Sound” on schedule2. The line in italics shows the change.

For the same example, the difference in **unsound** and **shallow** can be seen by looking at the caller **upgrade_prio** of **get_process** (that is syntactically unchanged) (either version in Figure 8). **unsound** flags a warning because it expects **get_process** to return different values in P_1 and P_2 for the `job` variable after its call to **get_process** (since the two versions of **get_process** are not equal), whereas **shallow** does not.

Finally, Figure 8 shows an example where a false warning was caused due to a missing specification of a callee (again from **schedule2**): whenever **get_process** returns a positive value in the `status` variable, the variable `job` is initialized to a non-null value. Hence, even with strong unsound assumptions made by **shallow** for the callees, modular DAC can still cause false warnings due to missing summaries.

Name	single	sound	unsound	shallow	nonmodular	versions	LOC	#procs
PrintTokens	5 (6)	5 (6)	0 (0)	0 (0)	0	5	565	18
PrintTokens2	6 (6)	3 (3)	0 (0)	0 (0)	0	10	508	19
Replace	32 (103)	10 (44)	4 (4)	4 (4)	2	32	562	21
Schedule	9 (17)	6 (14)	3 (3)	3 (3)	3	9	410	18
Schedule2	8 (16)	5 (36)	3 (7)	3 (3)	3	10	306	17
TotInfo	12 (12)	6 (8)	2 (2)	2 (2)	2	23	405	7
Space	38 (688)	15 (179)	10 (101)	10 (10)	MO	38	9128	136
Total	110 (848)	50 (290)	22(117)	22 (22)	10+	127	11884	236

Table 2: Name is the name of the benchmark; version is the number of different versions analyzed. LOC is lines of code and #procs is the number of procedures in each program. The numbers $x(y)$ mean that x versions and y procedures show warnings. “MO” is a out-of-memory exception.

Name	Diff	SymDiff	single	sound	unsound	shallow	nonmodular	LOC	#procs
firefly	1	1	1	1	1	1	1	634	7
moufilter	4	2	0	0	0	0	0	504	6
pciide	4	0	1	0	0	0	0	182	5
sfloppy	14	6	11	1	1	1	2	3404	20
diskperf	4	4	4	3	2	2	2	2319	24
event	1	1	0	0	0	0	1	555	5
cancel	3	1	0	1	0	0	0	476	5
Total	31	15	16	6	4	4	6	8074	72

Table 3: Name is the name of the benchmark; Diff is the number of procedures syntactically modified between Vista and Win7, SymDiff is the number of procedures for which the summary is true for the composed procedure. LOC is lines of code and #procs is the number of procedures in Win7 driver.

<pre> int upgrade_prio(prio, ratio) { int prio; float ratio; int status; struct process * job; if (prio < 1 prio > MAXLOPRIO) return (BADPRIO); if ((status = get_process (prio, ratio, &job)) <= 0) return (status); job->priority = prio + 1; ... } </pre>	<pre> int upgrade_prio(prio, ratio) { int prio; float ratio; int status; struct process * job; if ((status = get_process (prio, ratio, &job)) <= 0) return (status); job->priority = prio + 1; ... } </pre>
--	--

Figure 8: Imprecision in shallow

Table 3 shows the result of comparing two versions of sample device drivers in the Windows Device Driver Kit (WINDDK). The drivers for Windows Vista were considered as P_1 and the drivers for Windows 7 were considered as P_2 . The first column shows the name of the driver. The second column shows the number of procedures which were syntactically modified in going from Vista WINDDK to Win7 WINDDK for the same driver. The third column shows the number of functions which SYMDIFF failed to prove equivalent. Again the results are expected: the number of alarms are more for absolute correctness (**single**) than relative correctness. The sound strategy raises more alarms

than the unsound and shallow strategies. As mentioned earlier, the examples **sfloppy** and **event** illustrate that the set of warnings from **sound** does not always overapproximate **nonmodular**. For **cancel**, the option **sound** shows a warning whereas **single** does not — this happens due to the fact when mapped callees in the programs are called with different inputs, **sound** allows for the callee in the old program to pass and the callee in the new program to fail.

The experiments illustrate the feasibility of modular DAC towards providing a set of systematic knobs to narrow down the set of warnings resulting due to the program modification.

7. RELATED WORK

The idea of relative specifications is certainly not new; it goes back at least to checking simulation between two designs (usually at different levels of abstraction) using refinement mappings [1]. In contrast, DAC specifications are not necessarily refinement checks; the assertions present in a given program can be used to induce the relative specification. We use the reference program both to infer the unknown environment specification and also to help construct a modular proof. A concept similar to DAC has been explored in the context of filtering alarms for assertion checking of concurrent programs [17], however, this method applies only to bounded programs (see § 3 for details). Relative relaxed progress and memory safety have been formalized in the context of approximate program transformations [5, 6]; however little automation exists in checking them.

The most popular form of relative specifications for programs is equivalence checking. Such specifications come up

most naturally while checking for compiler optimizations using translation validation [26, 23, 18] and program refactoring [12, 25, 19] — however, such specifications are often too strong for most program changes during the course of evolution. Although DSE [25] provides differential summaries (for loop-free and recursion-free procedures) for arbitrary program changes, it does not provide a decision problem that DAC provides. In our experience with SYMDIFF, separating intended changes from unintended ones is the hardest problem when displaying differences to a user; DAC provides an intuitive specification whose violations are expected to be interesting for a user. Moreover, the DAC specifications need not be very program specific and can talk about relative specifications (such as using the *Valid* predicate for checking memory safety differentially) that are fairly abstract and thus applicable to most programs.

Product programs have been studied in the context of translation validation [32] and checking information flow properties [29]; these methods have been unified and generalized by recent works of Barthe et al. [4]. Similarly, several works on translation validation [26, 18, 33] infer simulation relations between synchronization points in two procedures to prove equivalence after intraprocedural transformations. However, these approaches only deal with intraprocedural transformations and do not account for interprocedural transformations. The construction of the composed program (§ 4) allows for specifying and inferring (intermediate) relative specifications for pairs of (possibly non-equivalent) procedures.

Finally, unlike previous approaches we provide a mechanism to leverage any off-the-shelf program verifier and invariant inference engine to check these relative specifications. The idea of comparing two programs with respect to assertions present has been suggested in previous works [20, 17, 5], but they do not provide a mechanism to specify or generate intermediate relative specifications, especially for loops and recursion. Mutual summaries [16] provide a mechanism for writing relative specifications by using quantified axioms to constrain the summaries of a pair of procedures. These mutual summaries can be seen as postconditions on the composed procedures. However, the approach cannot leverage off-the-shelf invariant inference engines to discover the intermediate relative specifications. On the other hand, [16] provides a modular checking for *relative termination* that is currently not handled by our DAC formulation. In the context of verifying safety of bug fixes, Gu et al. [14] investigate the *completeness* of a bug fix with distance-bounded weakest precondition, but cannot provide any soundness guarantees in the presence of unbounded loops and recursion.

8. CONCLUSION

In this work, we have described DAC as a mechanism for trading off cost for guarantees obtained while verifying evolving programs. We have reduced checking DAC to the analysis of a single program that can utilize standard program verification and invariant inference tools. We have provided an implementation of a simple scheme for automating the inference, and applied it towards verifying bug fixes and filtering alarms in real programs. We are currently integrating other tools based on interpolants [22] to generate relative specifications when the current scheme does not suffice.

APPENDIX

A. LOOPS

In this section, we describe how loops are transformed into tail-recursive procedures.² Although extracting loops as tail-recursive procedures is fairly standard, our approach differs from previous approaches [21] by avoiding the introduction of non-determinism in modeling the extracted procedure. This is important when comparing two programs; internal non-determinism makes program comparison difficult [17]. Our approach requires that the control flow graph is *reducible*, i.e., there is only one entry point for a loop. This assumption is true for almost any program generated from high-level languages such as C and Java.

We illustrate our approach informally using the example below where we use goto statements to model various control flow constructs present in high-level languages:

```
L0:
  s1;
  goto L0; //continue
  s2;
  goto L1; //break/jmp/return
  s3;
  goto L0; //loopback
L1:
  s4;
```

The loop is replaced by the following code fragment, where we use “[[s]]” to denote transforming any loops recursively inside a statement *s*.

```
L0: i' := call L0_loop(i);
  [[s1 :]]
  assume false; //goto L0;
  [[s2 :]]
  goto L1; //break/jmp/return
  [[s3 :]]
  assume false; //goto L0;
L1:
  [[s4 :]]
```

Here *i* represents the non-global variables in scope. In addition to the call to the tail recursive procedure `L0_loop`, the interesting aspect is the duplication of the last iteration of the loop body after the recursive call. The purpose of this is to handle goto statements that jump out of the loop (such as `goto L1`) [21]. The body of the tail recursive procedure transforms jumps to the loop head as tail-recursive calls. The main change to make the extracted procedure deterministic is to replace the jumps outside the loop by a statement that restores the state of the return and globals to the initial state.

```
proc L0_loop(i): i' {
  i' := i;
  [[s1 :]]
  i' := call L0_loop(i'); //tail-recursive call
  return;
  [[s2 :]]
  i' := i; g := old(g); return; //restore state
  [[s3 :]]
  i' := call L0_loop(i'); //tail-recursive call
  return;
}
```

²The exact Boogie options to be specified are `/print-Instrumented /extractLoops /deterministicExtract-Loops`.

B. REFERENCES

- [1] M. Abadi and L. Lamport. The existence of refinement mappings. *Theor. Comput. Sci.*, 82(2):253–284, 1991.
- [2] M. Barnett, B. E. Chang, R. DeLine, B. Jacobs, and K. R. M. Leino. Boogie: A modular reusable verifier for object-oriented programs. In *FMCO '05*, LNCS 4111, pages 364–387, 2005.
- [3] M. Barnett and K. R. M. Leino. Weakest-precondition of unstructured programs. In *PASTE '05*, pages 82–87, 2005.
- [4] G. Barthe, J. M. Crespo, and C. Kunz. Relational verification using product programs. In *Formal Methods (FM'11)*, LNCS 6664, pages 200–214, 2011.
- [5] M. Carbin, D. Kim, S. Misailovic, and M. C. Rinard. Proving acceptability properties of relaxed nondeterministic approximate programs. In *Programming Language Design and Implementation (PLDI'12)*, pages 169–180. ACM, 2012.
- [6] M. Carbin, D. Kim, S. Misailovic, and M. C. Rinard. Verified integrity properties for safe approximate program transformations. In *Partial Evaluation and Program Manipulation (PEPM'13)*, pages 63–66. ACM, 2013.
- [7] J. Condit, B. Hackett, S. K. Lahiri, and S. Qadeer. Unifying type checking and property checking for low-level code. In *POPL*, pages 302–314, 2009.
- [8] P. Cousot and R. Cousot. Abstract interpretation : A Unified Lattice Model for the Static Analysis of Programs by Construction or Approximation of Fixpoints. In *Symposium on Principles of Programming Languages (POPL '77)*. ACM Press, 1977.
- [9] L. de Moura and N. Bjørner. Z3: An efficient SMT solver. In *TACAS '08*, pages 337–340, 2008.
- [10] H. Do, S. G. Elbaum, and G. Rothermel. Supporting controlled experimentation with testing techniques: An infrastructure and its potential impact. *Empirical Software Engineering: An International Journal*, 10(4):405–435, 2005.
- [11] C. Flanagan and K. R. M. Leino. Houdini, an annotation assistant for esc/java. In *International Symposium of Formal Methods Europe (FME '01)*, 2001.
- [12] B. Godlin and O. Strichman. Regression verification. In *DAC*, pages 466–471, 2009.
- [13] S. Graf and H. Saïdi. Construction of abstract state graphs with PVS. In *Computer-Aided Verification (CAV '97)*, LNCS 1254, pages 72–83, June 1997.
- [14] Z. Gu, E. T. Barr, D. J. Hamilton, and Z. Su. Has the bug really been fixed? In *International Conference on Software Engineering (ICSE'10)*, pages 55–64. ACM, 2010.
- [15] B. Hackett, M. Das, D. Wang, and Z. Yang. Modular checking for buffer overflows in the large. In *International Conference on Software Engineering (ICSE'06)*, pages 232–241. ACM, 2006.
- [16] C. Hawblitzel, M. Kawaguchi, S. K. Lahiri, and H. Rebêlo. Towards modularly comparing programs using automated theorem provers. In *Conference on Automated Deduction (CADE'13)*, LNCS 7898, pages 282–299, 2013.
- [17] S. Joshi, S. K. Lahiri, and A. Lal. Underspecified harnesses and interleaved bugs. In *Principles of Programming Languages (POPL'12)*, pages 19–30. ACM, 2012.
- [18] S. Kundu, Z. Tatlock, and S. Lerner. Proving optimizations correct using parameterized program equivalence. In *Programming Language Design and Implementation (PLDI '09)*, pages 327–337. ACM, 2009.
- [19] S. K. Lahiri, C. Hawblitzel, M. Kawaguchi, and H. Rebêlo. Symdiff: A language-agnostic semantic diff tool for imperative programs. In *Computer Aided Verification (CAV'12)*, LNCS 7358, pages 712–717, 2012.
- [20] S. K. Lahiri, K. Vaswani, and C. A. R. Hoare. Differential static analysis: opportunities, applications, and challenges. In *Workshop on Future of Software Engineering Research (FoSER'10)*, pages 201–204. ACM, 2010.
- [21] A. Lal, S. Qadeer, and S. K. Lahiri. A solver for reachability modulo theories. In *Computer Aided Verification (CAV'12)*, LNCS 7358, pages 427–443, 2012.
- [22] K. L. McMillan. An interpolating theorem prover. In *Tools and Algorithms for the Construction and Analysis of Systems (TACAS'04)*, LNCS 2988, pages 16–30, 2004.
- [23] G. C. Necula. Translation validation for an optimizing compiler. In *ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '00)*, pages 83–94, 2000.
- [24] D. Notkin. Longitudinal program analysis. In *Workshop on Program Analysis For Software Tools and Engineering (PASTE '02)*, page 1. ACM, 2002.
- [25] S. Person, M. B. Dwyer, S. G. Elbaum, and C. S. Pasareanu. Differential symbolic execution. In *SIGSOFT FSE*, pages 226–237, 2008.
- [26] A. Pnueli, M. Siegel, and E. Singerman. Translation validation. In *Tools and Algorithms for Construction and Analysis of Systems (TACAS '98)*, pages 151–166, 1998.
- [27] Software-artifact Infrastructure Repository. Available at <http://sir.unl.edu/portal/index.html>.
- [28] O. Sokolsky and S. A. Smolka. Incremental model checking in the modal mu-calculus. In *Computer Aided Verification (CAV'94)*, LNCS 818, pages 351–363, 1994.
- [29] T. Terauchi and A. Aiken. Secure information flow as a safety problem. In *Static Analysis Symposium (SAS '05)*, LNCS 3672, pages 352–367, 2005.
- [30] The Boogie Verifier. Available at <http://boogie.codeplex.com>.
- [31] Verisec Suite. Available at http://se.cs.toronto.edu/index.php/Verisec_Suite.
- [32] A. Zaks and A. Pnueli. Covac: Compiler validation by program analysis of the cross-product. In *Formal Methods (FM'08)*, LNCS 5014, pages 35–51, 2008.
- [33] L. D. Zuck, A. Pnueli, B. Goldberg, C. W. Barrett, Y. Fang, and Y. Hu. Translation and run-time validation of loop transformations. *Formal Methods in System Design*, 27(3):335–360, 2005.