

# Inferring Annotations for Device Drivers from Verification Histories

Zvonimir Pavlinovic<sup>\*</sup>  
New York University, USA  
zvонimir@cs.nyu.edu

Akash Lal  
Microsoft Research, India  
akashl@microsoft.com

Rahul Sharma  
Stanford University, USA  
sharmar@cs.stanford.edu

## ABSTRACT

This paper studies and optimizes automated program verification. Detailed reasoning about software behavior is often facilitated by program invariants that hold across all program executions. Finding program invariants is in fact an essential step in automated program verification. Automatic discovery of precise invariants, however, can be very difficult in practice. The problem can be simplified if one has access to a candidate set of assertions (or *annotations*) and the search for invariants is limited over the space defined by these annotations. Then, the main challenge is to automatically generate quality program annotations.

We present an approach that infers program annotations automatically by leveraging the history of verifying related programs. Our algorithm extracts high-quality annotations from previous verification attempts, and then applies them for verifying new programs. We present a case study where we applied our algorithm to Microsoft's Static Driver Verifier (SDV). SDV is an industrial-strength tool for verification of Windows device drivers that uses manually-tuned heuristics for obtaining a set of annotations. Our technique inferred program annotations comparable in performance to the existing annotations used in SDV that were devised manually by human experts over years. Additionally, the inferred annotations together with the existing ones improved the performance of SDV overall, proving correct 47% of drivers more while running 22% faster in our experiments.

## CCS Concepts

•Theory of computation → Invariants; •Software and its engineering → Software verification; *Formal software verification*;

<sup>\*</sup> Author did part of the work as a research intern at Microsoft Research Bangalore, India

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. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [Permissions@acm.org](mailto:Permissions@acm.org).

ASE'16, September 3–7, 2016, Singapore, Singapore  
© 2016 ACM. 978-1-4503-3845-5/16/09...\$15.00  
<http://dx.doi.org/10.1145/2970276.2970305>

## Keywords

Program verification; Invariant generation; Learning invariants; Verification history; Big Code

## 1. INTRODUCTION

The performance of program verifiers depends on the discovery of precise assertions that hold during every program run, called program invariants. Examples of such assertions are procedure pre/post conditions and loop invariants. The task of finding invariants is often broken down into finding a set of candidate facts, or annotations, and then using these facts to establish invariants. For example, predicate-abstraction-based tools such as SLAM [2] or BLAST [13] rely on discovery of useful predicates to construct program invariants. Tools such as UFO [1] and Duality [20] rely on interpolation to generate candidates for procedure summaries. Each of these techniques, however, infer annotations by analyzing only the program given to be verified.

We propose a novel and complementary approach of inferring program annotations automatically by exploiting information available from prior verification runs. We build on the insight that annotations useful for verifying a particular program are often already observed earlier during the verification of related programs. For instance, programs that use the same API probably require similar annotations for verifying contracts of that API. We keep track of the verification history by accumulating a set of programs and the annotations required to construct their respective proofs. We leverage this history to generate a *small* set of annotations that are *useful* for subsequent (unseen) programs.

There are two key challenges in making this approach work. First, annotations are logical formulas over program variables, thus, tied to program-specific variable names. We abstract away from program-specific names by working with *abstract annotations*, which are arbitrary formulas with *holes*. Abstract annotations are *concretized* to a program by filling the holes with the program's variables.

The set of all (abstract) annotations in the verification history has the nice property that it is sufficient to establish the correctness of all programs observed in the history. However, this set is likely to be very large, making the verifier spend a significant amount of time just discarding invalid annotations. Our second challenge is to keep the set of inferred annotations small. We design a *minimization* algorithm that computes a set of abstract annotations such that: (1) it is enough to establish correctness proofs of all programs in the history, and (2) no smaller subset (or *syntactically simpler*

set, in a sense that we formalize later) is enough to establish all correctness proofs.

Our primary motivation behind these ideas is to improve the performance of Microsoft’s Static Driver Verifier (SDV) [21]. SDV is an industrial-strength tool for formal verification of Windows device drivers. SDV checks that drivers conform to certain properties (called *rules*) that establish correct usage of the Windows kernel APIs. SDV currently uses manually-tuned heuristics for obtaining a set of annotations that are passed to a program verifier. Using a repository of small in-house drivers, our techniques can not only replace the need for this manual effort, they even out-perform these heuristics and improve the performance of SDV overall.

We summarize our contributions as follows:

- Given a verification history, we formally define the notion of a minimal set of annotations, and present an algorithm for computing it.
- We apply the algorithm to SDV and experimentally show that inferring annotations from past verification efforts can potentially generate better annotations than ones provided by human experts. The set of abstract annotations inferred by our algorithm improved the verification times by 22% on average and reduced inconclusive results by 47% in our experiments.

## 2. OVERVIEW

This section motivates the need for automatic inference of program annotations and provides an overview of our techniques.

### 2.1 Static Driver Verifier

The Static Driver Verifier (SDV) has been an important success story for verification technology [2]. Over a decade, it has helped Windows developers to statically find bugs in Windows device drivers. The main verification engine of SDV was SLAM, with several upgrades along the way [2]. SDV switched to using an SMT-based verifier called Corral [15] for superior performance.

An overview of SDV is shown in Figure 1. SDV accepts the source code of a Windows device driver as input and links it against a model of the kernel (called “OS Model” in the figure). It then checks multiple rules that the driver must satisfy. These rules and kernel contracts (encoded in the OS Model) are made known to driver developers via Microsoft Developer Network (MSDN).<sup>1</sup> Each driver and rule results in (possibly multiple) programs with assertions, called *verification instances* that are fed to the verifier (Corral). Corral’s job is to find an execution that leads to an assertion violation in the verification instance. Such executions are reported as defects to the user.

Corral operates by lazily inlining procedures and utilizing an SMT solver to search through the partially-inlined program. To help Corral, SDV uses annotation-based invariant generation. SDV generates annotations based on the rule being checked and runs Houdini [9] on the verification instance to compute invariants constructed from these annotations. These invariants are injected back to the verification

instance as *assumptions*, which help Corral prune search without compromising soundness.

Corral has four possible outcomes, as mentioned in Figure 1. Corral uses over-approximations (refined by invariants inferred by Houdini), hence it can prove correctness and return “proof”. Corral can also find a bug and report the failing execution. Corral stops search when it hits an internal coverage bound [18] and returns “bound”. Although this is an inconclusive verdict, it is still considered more useful than a Timeout outcome because the latter does not guarantee any coverage. SDV has experienced several upgrades, including major improvements to Corral [16, 17]. Showing further improvements would truly build on the state-of-the-art in the area.

This paper focuses on improving annotation generation in SDV. The current process of generating annotations is guided by a set of heuristics that have been manually tuned and maintained over several years. These heuristics can be found in our technical report [28]. What makes this approach feasible is that the heuristics only look at the rule being checked (not the driver) for generating the annotations. While driver code is not known ahead of time, the set of rules is fixed and known to SDV developers before SDV’s release.

The goal of this paper is two-fold. First, we want to automatically generate annotations instead of requiring manual effort. Second, we wish to show that automated inference can out-perform the expert-driven heuristics. We achieve these goals by learning useful information from past verification efforts.

### 2.2 Examples of SDV Rules

We give two examples of rules (properties) that SDV checks. The rules are described in an abstract manner, without using the actual tool notation, for clarity.

**Acquire-release rules.** The OS kernel provides multiple resources to help the driver accomplish its task. The kernel expects that the driver (specifically, the dispatch routines of the drivers) must release all acquired resources when it exits. SDV checks this property by instrumenting the Acquire-Release API as shown in Figure 2(a). The actual code of the API is immaterial while checking this property, thus, is not shown. The rule introduces the *model* variable depth to keep track of the number of resources held; the variable is asserted to be zero when the driver’s dispatch routine exits. This rule abstractly captures the *SpinLockRelease* rule of SDV.<sup>2</sup>

Figure 2(b) shows a family of (fake) drivers that exercise this API correctly. Entry points of the drivers are the procedures  $P_n$  and  $Q$ .  $P_n$  is parameterized by the value of  $n$ . It calls the routine  $\text{dispatch}_{P_n}$  where we use the notation  $[st]^n$  to denote  $n$  occurrences of the statement  $st$ . These programs model typical usage of the Acquire-Release API that we have seen in drivers (while abstracting away irrelevant details). SDV links the driver code (Figure 2(b)) together with the instrumented API code (Figure 2(a)) to produce a verification instance that is fed to Corral. The construction of a proof of correctness and the role of annotations is discussed in the next section.

**Irql-based rules.** An *Interrupt Request Level* (IRQL) captures the priority associated with a task. Tasks with higher IRQL cannot be interrupted by tasks with a lower

<sup>1</sup> [https://msdn.microsoft.com/en-us/library/windows/hardware/ff552840\(v=vs.85\).aspx](https://msdn.microsoft.com/en-us/library/windows/hardware/ff552840(v=vs.85).aspx)

<sup>2</sup> [https://msdn.microsoft.com/en-us/library/windows/hardware/ff552780\(v=vs.85\).aspx](https://msdn.microsoft.com/en-us/library/windows/hardware/ff552780(v=vs.85).aspx)

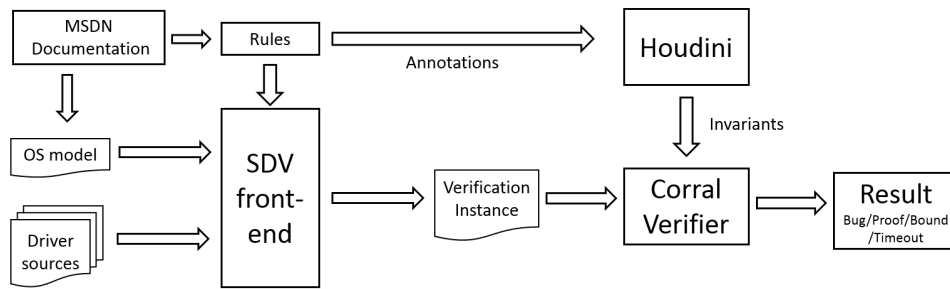


Figure 1: An overview of SDV’s implementation.

IRQL. Drivers often raise the IRQL level to perform critical activity uninterrupted, but are required not to spend too much time at the high IRQL for the sake of responsiveness of the system. SDV checks that drivers do not call certain time consuming kernel APIs when at a high IRQL<sup>3</sup>.

Figure 3 shows IRQL modeling and usage. The OS model variable `sdv_irql_current` records the current IRQL value of the processor. Kernel APIs `KeRaiseIrql`<sup>4</sup> and `KeLowerIrql` can be used to change the IRQL value. The rule simply asserts that the procedure “do\_work\_at\_low\_irql” is only called when IRQL is 0 (lowest possible value). The program with entry point `main` is a (fake) driver that correctly exercises the kernel API.

### 2.3 Annotations and Proofs

We now set up some notation to describe program proofs. If  $f$  is a procedure and  $\phi$  is a formula, we use  $[\phi]@f$  to denote that  $\phi$  is a valid postcondition of  $f$ . If  $f$  has a loop starting at location  $L$  then the notation  $[\phi]@f@L$  denotes that  $\phi$  is a valid loop invariant for  $L$ . We model assertion failures as setting of a special `ok` bit to `false`. For a variable  $x$ , let `old(x)` refer to the value of  $x$  at the beginning of the procedure or loop, depending on the context in which it is used. For instance,  $[(x == \text{old}(x) + 1) \vee \neg \text{ok}]@f$  means that the execution of  $f$  either increments the value of  $x$  or it fails an assertion. In other words, if  $f$  doesn’t fail then it increments  $x$ .  $[x == \text{old}(x)]@f@L$  means that the loop at location  $L$  of procedure  $f$  preserves the value of  $x$  across an arbitrary number of loop iterations. A *proof* of correctness of a program is simply a sequence of mutually-inductive postconditions of procedures or loops in the program that imply `ok == true` at the end of the program. For simplicity (and without loss of generality) we do not talk about procedure preconditions in this paper. *Annotations* are simply formulas that serve as candidates for postconditions.

Figure 2(c) shows possible proofs for the previously introduced programs of Figure 2. The figure contains two possible proofs for  $P_n$  (marked as “A” and “B”) and a single proof for  $Q$ . Note that all postconditions in proof A of  $P_n$  do not depend on the value of  $n$ , whereas proof B is specific to the value of  $n$ .

Generation of invariants from a given set of annotations (which we call *annotation-based invariant generation*) is much

simpler than full-blown verification, often even decidable. One may use, for example, predicate abstraction [3] to construct invariants that are Boolean combinations of the given annotations. In our work, we use the Houdini algorithm [9] to find conjunctive invariants: ones that are conjunctions of some subset of the given annotations. This problem has a lower complexity than predicate abstraction and is very fast in practice for small to medium number of annotations. For example, given annotations  $\{\text{depth} == 0, \psi_{-1}, \psi_0, \psi_1, \eta\}$  (not knowing if they are valid postconditions or loop invariants) it is very efficient to *construct* a proof for  $P_n$  (of type A) using Houdini.

### 2.4 Minimal Repositories

Our technique requires a repository of programs and their proof of correctness. The proofs may be constructed manually or by using proof-generating verifiers. We do not expect to control the proof-generation process. Suppose we have programs  $Q$  and  $P_n$  for each  $n \in N$ , for some large set  $N$ . Further, suppose we have proof of type B (Figure 2(c)) for  $P_n$  for all values of  $n$ , except  $n_0$ , and  $P_{n_0}$  has a proof of type A. These  $N$  programs together with their corresponding proofs constitute a repository.

We extract all annotations present in the proofs, which produces a large set  $\mathcal{A} = \{\text{depth} == 0, \psi_{-1}, \psi_0, \psi_1, \eta, \psi_0 \wedge \eta\} \cup \{\sigma_n \mid n \in N - \{n_0\}\}$ . Retaining a large set of annotations is inefficient, even for annotation-based invariant generation techniques. Moreover, some of the annotations are very specific to a program, e.g.,  $\sigma_n$  is only useful for proving correctness of  $P_n$ .

Our technique minimizes  $\mathcal{A}$  while retaining its invariant-generation power. The “power” is captured using a *cost* metric based on the ability of a set of annotations to prove a set of programs correct, given a fixed verifier. The cost is  $\infty$  if some program cannot be proved, otherwise, it reflects the running time of the verifier. Our algorithm simplifies  $\mathcal{A}$  by dropping annotations or making them syntactically simpler as long as the cost does not increase (or only increases by a *tolerable* amount; the exact formulation can be found in Section 3).

For illustration, assume that the cost becomes  $\infty$  as soon as the annotations cannot establish some loop invariant or postcondition of a recursive procedure (intuitively, because these are the critical parts of a proof), and is unit cost otherwise. Starting with  $\mathcal{A}$ , our algorithm drops `depth == 0` and  $\eta$  from this set because these are not important for the inductive argument; i.e., cost remains unit after dropping them. Next, if it tries to drop  $\psi_0$ , the cost becomes  $\infty$  be-

<sup>3</sup>See, for example, [https://msdn.microsoft.com/en-us/library/windows/hardware/ff547747\(v=vs.85\).aspx](https://msdn.microsoft.com/en-us/library/windows/hardware/ff547747(v=vs.85).aspx)

<sup>4</sup>[https://msdn.microsoft.com/en-us/library/windows/hardware/ff553079\(v=vs.85\).aspx](https://msdn.microsoft.com/en-us/library/windows/hardware/ff553079(v=vs.85).aspx)

```

var depth: int;

procedure init() {
  depth := 0;
}

procedure Acquire() {
  depth := depth + 1;
}

procedure Release() {
  depth := depth - 1;
}

procedure d_exit()
{
  assert depth == 0;
}

(a)

// Definitions
 $\psi_i \equiv \text{depth} - \text{old}(\text{depth}) == i$ 
 $\sigma_n \equiv \text{old}(\text{depth}) == n \Rightarrow \text{depth} == n$ 
 $\eta \equiv (\text{old}(\text{depth}) == 0 \Rightarrow \text{ok})$ 

// Proof A of Pn
[depth == 0]@init, [ $\psi_1$ ]@Acquire, [ $\psi_{-1}$ ]@Release, [ $\eta$ ]@d_exit, [ $\psi_0$ ]@dispatchPn@L1, [ $\eta$ ]@dispatchPn

// Proof B of Pn
[depth == 0]@init, [ $\psi_1$ ]@Acquire, [ $\psi_{-1}$ ]@Release, [ $\eta$ ]@d_exit, [ $\sigma_n$ ]@dispatchPn@L1, [ $\eta$ ]@dispatchPn

// Proof of Q
[depth == 0]@init, [ $\psi_1$ ]@Acquire, [ $\psi_{-1}$ ]@Release, [ $\eta$ ]@d_exit, [ $\psi_0 \wedge \eta$ ]@dispatchQ@L2, [ $\psi_0 \wedge \eta$ ]@dispatchQ

(c)

```

```

procedure Pn()
{ call init(); call dispatchPn(); }

procedure dispatchPn() {
  [call Acquire()]n;
L1: while (*)
  { call Acquire(); call Release(); }
  [call Release()]n; call d_exit();
}

procedure Q()
{ call init(); call dispatchQ(); }

procedure dispatchQ() {
L2: while (*)
  { call Acquire(); call Release();
    call dispatchQ(); }
  call d_exit();
}

(b)

```

**Figure 2: (a) An Acquire-Release API, (b) a family of programs exercising the API, and (c) possible proofs of correctness of the programs.**

cause the loop invariant of  $P_{n_0}$  is lost. Thus,  $\psi_0$  is retained in  $\mathcal{A}$ . Next, each of the  $\sigma_n$  annotations get dropped. Even though these annotations were loop invariants in the original proofs, they can be replaced by the more general annotation  $\psi_0$  that is present in  $\mathcal{A}$ . In this way, learning from a large set of proofs increases the chances of finding annotations useful for many programs.

Finally, while  $\psi_0 \wedge \eta$  cannot be dropped, our algorithm tries to simplify its Boolean structure. The algorithm simplifies it to  $\eta$  because having annotations  $\{\psi_0, \eta\}$  is enough for annotation-based invariant generation to establish  $\psi_0 \wedge \eta$  as an invariant. At this point, the algorithm reaches a fixpoint where no annotation can be dropped or simplified and it returns the set:  $\{\psi_{-1}, \psi_0, \psi_1, \eta\}$ .

The algorithm is non-deterministic; it could have chosen to drop  $\psi_0 \wedge \eta$  in its first iteration because  $\eta$  was still present in  $\mathcal{A}$ . In general, our algorithm only guarantees a locally optimal solution with respect to a given cost metric. Globally-optimal solutions are also possible to compute, but at a higher cost, which was not justified in our experiments.

Although the programs considered here are simple, they are derived from real-world code. They reflect common usage patterns of acquire-release-kind of APIs that we have observed in drivers. The program  $Q$  illustrates an uncommon (but not rare) scenario where recursion happens via kernel callbacks (driver code itself typically does not exhibit recursion).

## 2.5 Abstract Annotations

Annotations are formulas over program variables. In general, different programs have different variables. To abstract away from program-specific variables, we introduce the concept of an *abstract annotation* that is a formula over only *generic* and *shared* variables.

We call the set of global variables common to all programs in the repository as the *shared vocabulary*. We assume these variables serve a similar role in all programs, e.g., the `depth` variable of Figure 2 will be present in all programs that exercise the acquire-release API, and `sdv_irq1_current` will be present in all programs that exercise the IRQ1 API of Figure 3. Shared variables, i.e., variables in the shared

```

// DRIVER
procedure main() {
  var loc: int;
  call init();
  call loc := KeRaiseIrql(2);
  while(*) {
    call KeLowerIrql(loc);
    call do_work_at_low_irql();
    call loc := KeRaiseIrql(2);
  }
  call KeLowerIrql(loc);
}

// RULE
procedure do_work_at_low_irql()
{ assert sdv_irql_current == 0; }

// OS MODEL and RULE
var sdv_irql_current: int;

procedure init()
{ sdv_irql_current := 0; }

procedure KeRaiseIrql(new_irql: int)
returns (old_irql: int) {
  old_irql := sdv_irql_current;
  sdv_irql_current := new_irql;
}

procedure KeLowerIrql(new_irql: int)
{ sdv_irql_current := new_irql; }

```

**Figure 3: IRQL modeling and example of usage.**

vocabulary, can be freely used in annotations because they are present in all programs that exercise the same rule.

Generic variables are not specific to any program. There are four kinds of generic variables:

{LOCAL, GLOBAL, FORMALIN, FORMALOUT}.

The kind of a generic variable is determined statically by inspecting the declaration of a program variable. An annotation is converted to an abstract annotation by replacing variables by generic variables of the corresponding type. For example, consider the postcondition  $[x == y]@f$  on a procedure  $f$  with formal input argument  $x$  and formal output argument  $y$ . This will get converted to the abstract annotation  $(\$fin == \$fout)$  where  $\$fin$  and  $\$fout$  are generic variables of kind FORMALIN and FORMALOUT, respectively.

Abstract annotations are concretized when applied to a program. Let  $p$  be a program and  $proc$  a procedure in  $p$ . We define a concretization function  $\gamma_{p,proc}$  as follows. For an abstract annotation  $a$ ,  $\gamma_{p,proc}(a)$  returns all annotations such that a generic variable of type GLOBAL is substituted with some global variable of  $p$ , a generic variable of type FORMALIN is substituted with some formal-in parameter of  $proc$ , and similarly for FORMALOUT and LOCAL.  $\gamma_{p,proc}$  must return all such annotations.  $\gamma_{p,proc}$  leaves shared variables unchanged.

Consider the driver in Figure 3. Proving correctness of `main` requires a loop invariant that the value of

`sdv_irql_current` is unchanged across loop iterations, which in turn requires that the value of `loc` is maintained across iterations. That is, it requires the loop invariant  $loc == old(loc)$ . Our technique, once it observes a proof with this loop invariant will generate and keep the abstract annotation  $\$floc == old(\$floc)$ , where  $\$floc$  is a generic variable of type LOCAL. With this abstract annotation, when SDV is executed on a driver that exercises the IRQL rule, the annotation will get instantiated with all local variables manipulated by loops in the driver, and the annotation-based invariant generation algorithm will be able to establish correctness of drivers that require preservation of local variable values, similar to `main` of Figure 3.

### 3. ALGORITHM

We now formally describe the annotation inference algorithm. Our presentation of the algorithm is general, abstracting away from SDV specifics for ease of presentation and to emphasize the use of verification histories. However, we evaluate the algorithm only on SDV in this work. We leave generalization to other verification domains as future work. We start by introducing the necessary notation and definitions.

**Language.** We assume an imperative programming language with standard features such as global variables, procedures, `assume` and `assert` statements, assignments, etc. We also assume that programs in this language do not have loops. Loops can be encoded using recursion. This allows our framework to only concentrate on procedure postconditions for establishing proofs of correctness.

Given a program  $p$ , we denote the set of procedures in  $p$  with  $procs(p)$ . Each procedure can be *annotated* with any number of first order logic (FOL) formulas. These formulas are defined over procedure parameters and global variables and they do not take part in program execution; they are used by program verifiers as candidate postconditions for establishing program correctness.

**Abstract annotations.** Let  $V$  be a set of variables called the shared vocabulary. All programs must contain  $V$  as global variables. Let  $G$  be a set of generic variables. None of the programs contain a variable from  $G$ . An *abstract annotation*  $a$  is a formula over variables in  $V \cup G$ . Further, for every program  $p$  and procedure  $proc \in procs(p)$ , we assume a function  $\gamma_{p,proc}$  that maps an abstract annotation to a set of concrete annotations. As defined in the previous section,  $\gamma_{p,proc}$  substitutes generic variables with the variables in scope of  $proc$ .

We call a finite set of abstract annotations  $t \in \mathfrak{T}$  a *template*. Given a program  $p$  and a template  $t$ ,  $annotate(p, t)$  returns  $p$  where each procedure  $proc \in procs(p)$  is annotated with the set of program annotations  $\bigcup_{a \in t} \gamma_{p,proc}(a)$ .

**Verification.** Given a fixed verifier, for a program  $p$  and a template  $t$ , we say that  $proves(p, t)$  holds if the verifier can prove the correctness of  $annotate(p, t)$ . The verifier can use the procedure annotations as potential postconditions during the verification. Also, we require that if  $proves(p, t)$  and  $t \subseteq t'$ , then  $proves(p, t')$  must hold as well. In other words, if a set of abstract annotations is sufficient for proving some program correct then all of its supersets are also sufficient.

#### 3.1 Problem

Our annotation inference problem is defined using the notion of an objective relation. Such relations are used to

encode what templates (and hence annotations) are more desirable for the current application in mind.

**DEFINITION 1** (OBJECTIVE RELATION FOR A PROGRAM).  
*We say  $\rightarrow_p: \mathfrak{T} \times \mathfrak{T}$  is an objective relation for a program  $p$  iff (1) it is well-founded and (2) for each  $t \in \text{dom}(\rightarrow_p)$ , the set  $\{t' \mid t \rightarrow_p t'\}$  is finite.*

The objective relation is hence *finite branching*. One example of such a relation is the proper subset relation, i.e.,  $t \rightarrow_p t'$  iff  $t' \subset t$ . Another example would be the relation where  $t'$  is a copy of  $t$  except that an annotation in  $t'$  is a sub-formula of the corresponding annotation in  $t$ . This relation roughly corresponds to the *syntactically simpler* concept mentioned in Section 1. Section 4 presents the objective relation used in our experiments with SDV.

We extend the objective relation to a set  $P$  of programs  $p_1, \dots, p_n$ :

**DEFINITION 2** (OBJECTIVE RELATION FOR PROGRAMS).  
*We define  $\rightarrow_P: \mathfrak{T} \times \mathfrak{T}$ , an objective relation for a set of programs  $P = \{p_1, \dots, p_n\}$  as a well founded and finite branching relation  $t \rightarrow_P t' \Leftrightarrow \bigwedge_i t \rightarrow_{p_i} t'$ .*

We proceed by defining the notion of a minimal template that intuitively stands for a locally optimal template. The locality is defined as a branching set of a template induced by a given objective relation.

**DEFINITION 3** (MINIMAL PROGRAM TEMPLATE).  
*Given a program  $p$ , a template  $t$  such that  $\text{proves}(p, t)$ , and an objective relation  $\rightarrow_p$ , we say  $t'$  is a minimal template iff:*

1.  $\text{proves}(p, t')$
2. there exists no  $t''$  such that  $t' \rightarrow_p t''$  and  $\text{proves}(p, t'')$

The definition states that a minimal template must prove a given program and none of its immediate ( $\rightarrow_p$ ) successors do. We point out that, in the above definition,  $t$  only ensures that  $p$  is correct and this definition establishes no relationship between a minimal template and  $t$ . However, the results computed by the implementations for finding minimal templates can be dependent on  $t$ .

Our inference algorithm is built around the notion of a minimal template. We hence define the problem of finding a minimal template for a given program.

**PROBLEM 1** (COMPUTING A MINIMAL TEMPLATE).  
*Given a program  $p$ , a template  $t$  such that  $\text{proves}(p, t)$ , and an objective relation  $\rightarrow_p$ , the problem of computing a minimal template is finding a formula  $t'$  that is minimal subject to  $p$  and the ordering  $\rightarrow_p$ .*

We define a *program repository* as  $R = [(p_1, t_1), \dots, (p_n, t_n)]$  where  $\text{proves}(p_i, t_i)$  for  $1 \leq i \leq n$ . Repositories capture verification histories. The set of programs in the repository  $R$  is denoted by  $P_R$ . The actual technique used for proving the correctness of  $p_i$  can be arbitrary. However, we envision that in practice, a verifier will be fixed for the whole repository. We now define a locally optimal template for a verification history.

**DEFINITION 4** (MINIMAL REPOSITORY TEMPLATE).  
*Given a program repository  $R = [(p_1, t_1), \dots, (p_n, t_n)]$  where  $\text{proves}(p_i, t_i)$  for all  $1 \leq i \leq n$ , and an objective relation  $\rightarrow_{P_R}$ , we say that  $T$  is a minimal repository template (subject to  $\rightarrow_{P_R}$ ) iff the following holds*

---

**Algorithm 1** Computing a minimal program template

---

**Require:**  $\text{proves}(p, t)$  and  $\rightarrow$  is an objective relation

```

1: procedure MINTEMPLATE( $p, t, \rightarrow$ )
2:    $mint \leftarrow t$ 
3:   loop
4:     for all  $t' \in \{t' \mid t \rightarrow t'\}$  do
5:       if  $\text{proves}(p, t')$  then
6:          $mint \leftarrow t'$ 
7:       goto 3
8:   return  $mint$ 

```

---

1.  $\text{proves}(p_i, T)$  for all  $1 \leq i \leq n$
2. there exists no  $T'$  such that  $T \rightarrow_{P_R} T'$  and  $\text{proves}(p_i, T')$  for all  $1 \leq i \leq n$

We are now ready to formally state the problem of inferring program annotations from past verification runs.

**PROBLEM 2** (INFERRING PROGRAM ANNOTATIONS).  
*Given a program repository  $R = [(p_1, t_1), \dots, (p_n, t_n)]$  where  $\text{proves}(p_i, t_i)$  for all  $1 \leq i \leq n$ , and an objective relation  $\rightarrow_{P_R}$ , the problem of inferring program annotations is to find a template  $T$  that is a minimal repository template subject to  $\rightarrow_{P_R}$ .*

In the sequel, we suppress the subscripts of the  $\rightarrow$  relations for brevity. We now show algorithms for solving the problems of computing minimal templates.

## 3.2 Solution

We start with the solution for Problem 1 shown in Algorithm 1. The `MinTemplate` algorithm assumes that a given template  $t$  is sufficient to establish correctness of  $p$  and that  $\rightarrow$  is an objective relation. We start by considering  $t$  as a minimal template candidate (line 2). We continue by enumerating all immediate successors of  $t$  by  $\rightarrow$  (line 4). Then, the algorithm checks if any of the successors can prove  $p$  (line 5). If so, then the algorithm sets such a successor as a candidate for the minimal template and repeats the whole process (lines 6 and 7). Otherwise, the minimal candidate is returned as the solution (line 8).

**THEOREM 1.** *Let  $p$  be a program,  $t$  a template, and  $\rightarrow$  an objective relation. If  $\text{proves}(p, t)$ , then Algorithm 1 computes a minimal template for  $p, t$ , and  $\rightarrow$ .*

**PROOF.** Since  $\rightarrow$  is finite branching, we have that inner loop at line 4 terminates. From the fact that  $\rightarrow$  is well-founded, it follows that the outer loop at line 3 also terminates. Since lines 4 and 5 simply follow the definition of a minimal template, we have that the returned template is indeed minimal.  $\square$

The complexity of the algorithm depends on  $\rightarrow$  relation and implementation of  $\text{proves}$ . Assuming  $\text{proves}$  has unit complexity, the running time of Algorithm 1 is  $O(l \cdot m)$ , where  $l$  is the longest well-founded chain of  $\rightarrow$  and  $m$  is the maximum size of the branching sets  $\max\{|\{t' \mid t \rightarrow t'\}| \mid t \in \text{dom}(\rightarrow)\}$ . However, proving a program correct is undecidable in general and expensive in practice. Further, annotating a program  $p$  given a template  $t$  can also be expensive if for  $\text{proc} \in \text{procs}(p)$ ,

---

**Algorithm 2** Computing a minimal repository template

---

**Require:**  $proves(p_i, t_i)$  for all  $(p_i, t_i) \in R$

**Require:**  $\rightarrow$  is an objective relation

```
1: procedure MINREPOTEMPLATE( $R, \rightarrow$ )
2:    $mints \leftarrow []$ 
3:   for all  $i \in [1, \dots, |R|]$  do
4:      $(p, t) = R[i]$ 
5:      $mints[i] = MinTemplate(p, t, \rightarrow)$ 
6:    $C = \emptyset$ 
7:   for all  $i \in [1, \dots, |R|]$  do
8:      $C = C \cup mints[i]$ 
9:     for all  $t \in \{t' \mid C \rightarrow_{\{p_1, \dots, p_i\}} t'\}$  do
10:       $b = true$ 
11:      for all  $j \in [1, \dots, i]$  do
12:         $b = b \wedge proves(p_j, t)$ 
13:      if  $b$  then
14:         $C = t$ 
15:      goto 9
16:   return  $C$ 
```

---

the concretization function  $\gamma_{p,proc}$  has a large image;  $t$  can then potentially be instantiated with a large number of concretizations. This high complexity of the algorithm can be remedied in practice by choosing  $\gamma_{p,proc}$  with smaller images and exploiting the structure of  $\rightarrow$  if the relation is known beforehand. We note that a result computed by the algorithm is not necessarily minimum. As we show in Section 5, minimal templates sufficed for all practical purposes in our experiments.

Algorithm 2 computes a minimal repository template by building on Algorithm 1. First, we find the minimum template for each program in the repository and store it in  $mints$  (lines 3-5). Next, the minimal repository template is set to the empty set of clauses (line 6). The outer loop (lines 7-15) has the invariant that after the  $i^{th}$  iteration  $C$  is a minimal repository template for the sub-repository  $[(p_1, t_1), \dots, (p_i, t_i)]$ . The inner loop checks if an immediate successor of  $C$  can prove the correctness of the programs  $p_1, \dots, p_i$ . If so, then  $C$  is updated to that successor template.

One possible optimization is to cache the clauses under which a program can/cannot be proved. Therefore, if  $proves(p, t)$  is in the cache and we later make a query  $proves(p, t')$  where  $t \subseteq t'$ , then, we can return *true*. Similarly, if  $\neg proves(p, t)$  is in the cache and we make a query  $proves(p, t')$  where  $t' \subseteq t$ , then, we can return *false*.

**THEOREM 2.** *Let  $R$  be a program repository and  $\rightarrow$  an objective relation. If  $proves(p, t)$  for all  $(p, t) \in R$ , then Algorithm 2 computes a minimal repository template for  $R$  and  $\rightarrow$ .*

**PROOF.** For every subrepository  $R_i = [(p_1, t_1), \dots, (p_i, t_i)]$  the relation  $\rightarrow_{R_i}$  is well-founded and finite branching for any chosen reduction operator. Since the body of the loop at line 9 follows the definition of a minimal repository template, then at the end of each iteration of the loop at line 7,  $C$  is a minimal template for  $R_i$ , as pointed out earlier. The result then follows from the case when  $i = |R|$ .  $\square$

We also point out that the loop starting at line 7 is in fact not necessary for optimality. The algorithm can immediately start with  $C$  as the union of all minimal program templates stored in  $mints$ . However, such a  $C$  could become impractically large. In Algorithm 2, the size of  $C$  is kept moderate. Observe that the minimal templates computed using these two versions of the algorithm might not be the same. This is because minimal templates are not unique.

## 4. APPLICATION TO SDV

We use a specific objective relation to infer useful program annotations for SDV on Windows device drivers.

### 4.1 Repository

For internal testing, SDV uses a set of toy drivers, collectively called the *Rule Test Suite* (RTS) for quick “smoke testing” of SDV. These drivers are small, often a few hundred lines of code (with relevant part in tens of lines of code only). We use RTS as the program repository for inferring a minimal repository template for each rule. We will refer to such templates as “inferred templates” and to manually crafted templates simply as “manual templates”. The use of RTS as the training set is quite natural as it allows us to leverage existing test cases to learn and improve performance on real drivers. RTS consists of 304 drivers totaling around 100KLOC.

Once we infer a template, we measure the performance of SDV using the template on a set of real device drivers. We use 66 device drivers for this task, totaling around 700KLOC. We note that the device drivers were functionally very different: they included storage drivers, modem drivers, bus drivers, etc.

SDV has hundreds of rules. For this paper, we concentrate on a collection of 28 rules. Of these, 14 rules were selected because they were known to cause performance issues for SDV. We added, randomly chosen, 14 other rules. These 28 rules on the 66 drivers produced a total of 1420 verification instances (not all rules apply to all drivers).

**Obtaining Annotations.** The RTS suite consists of fairly small drivers that are easy to prove manually. We, however, automated the entire process by running a proof-generating verifier. Our implementation uses Duality [20] but conceptually we could have used tools such as SLAM and Yogi [4] as well.

Our objective is to learn a template per rule. For each of the 28 rules, we (1) form a repository of programs that assert the rule and their corresponding correctness proofs and (2) define the shared vocabulary  $V$  to consist of model variables of the rule as well as OS model variables. Every program in a repository is guaranteed to include the corresponding shared vocabulary as global variables.

### 4.2 Objective Relation

We restrict abstract annotations to be *clauses*, i.e., a disjunction of formulas. Conjunctions at the top level are broken down into multiple annotations, one for each conjunct. For convenience, we think of a clause as a set of formulas where disjunction is implied between the elements of the set. The empty set corresponds to `true`. Given two annotations  $a_1$  and  $a_2$ , by  $a_1 \subseteq a_2$  we therefore designate that the formulas of  $a_1$  are a subset of the formulas of  $a_2$ . For a template  $t$ , we define  $\bar{t}$  to be  $t$  consistently indexed by the set  $\{1, \dots, |t|\}$ . In other words,  $\bar{t} = \{a_1, \dots, a_{|t|}\}$  where  $a_i \in t \iff a_i \in \bar{t}$ .

For convenience, we define  $\bar{t}[i] = a_i$ . Given two templates  $t_1$  and  $t_2$ , we say  $t_2$  is *simplified* (or *simpler*) than  $t_1$ , written  $t_2 \preceq t_1$ , iff (1)  $|t_1| = |t_2|$  and (2)  $\bar{t}_2[i] \subseteq \bar{t}_1[i]$  for all  $1 \leq i \leq n$ . If  $t_2 \preceq t_1$  and there exists  $i$  such that  $\bar{t}_2[i] \subset \bar{t}_1[i]$ , we say  $t_2$  is *strictly simpler* than  $t_1$ , written  $t_2 \triangleleft t_1$ . Finally,  $t_2 \triangleleft_1 t_1$  holds iff  $t_2 \triangleleft t_1$  and  $|\bar{t}_1/\bar{t}_2| = 1$ . In other words,  $t_2$  is strictly simpler than  $t_1$  but only at one abstract annotation; we then say  $t_2$  is 1-simpler than  $t_1$ . For example, suppose a template  $t'$  is  $t$  except that literal  $l$  is in  $\bar{t}[1]$  but not in  $\bar{t}'[1]$ . Then we have  $t' \preceq t$ ,  $t' \triangleleft t$ , and  $t' \triangleleft_1 t$ .

Given a program  $p$  and a template  $t$ , let  $h(p, t)$  be the result of running Houdini on  $\text{annotates}(p, t)$ . Further, let  $c_{\text{perf}}(p, t)$  denote the number of procedures that Corral inlined if it was able to prove correctness of  $h(p, t)$ , and  $-1$  otherwise.  $c_{\text{perf}}(p, t)$  measures the performance of Corral. We use it as a proxy for the running time that is independent of the machine configuration. We now define the *corral* objective relation. Given a program  $p$  and two templates  $t_1$  and  $t_2$ ,  $t_1 \rightarrow_p t_2$  holds iff:

- $t_2 \triangleleft_1 t_1$
- $0 \leq c_{\text{perf}}(p, t_2) \leq 2 * c_{\text{perf}}(p, t_1)$

The above definitions encode our intention to find those templates that are structurally simpler and smaller, if one views making a clause empty as removing it. The reason why we chose 1-simpler relation instead of the general simplification is that we want to keep branching sets of the objective relation tractable. This way, we are sacrificing optimality for better inference times. Also note that our objective relation allows the performance of Corral to get worse when using  $t_2$ . But Corral must still be able to prove correctness. We allow the degradation in performance to allow more opportunities for the templates to get simplified. However, we still restrict the performance to not get out of hand (not more than a factor of 2). This allows us to kill the execution of Corral on  $h(p, t_2)$  as soon as it inlines twice as many procedures as on  $h(p, t_1)$ , without waiting for a verdict.

Lastly, we define a repository objective relation. For a set of programs  $P = \{p_1, \dots, p_n\}$  we define  $t_1 \rightarrow_P t_2 \iff \bigwedge_i t_1 \rightarrow_{p_i} t_2$ . In other words, we want to infer templates that are consistently optimal for all repository programs, subject to the objective relation  $\rightarrow_p$ .

### 4.3 Implementation

We use a slightly modified version of Algorithm 2 for computing a minimal repository template. We describe the algorithm by only explaining the modifications that we introduced.

We first note that the *corral* objective relation has to be computed by actually running Corral. While performing a greedy descent, we simply enumerate all 1-simplified templates and run Corral on each of them. This run tells us whether the program can be proved and the number of inlined procedures, if any. Also, we enumerate first those annotations where the corresponding simplified clause is empty. Then, we enumerate annotations where the simplified clause has one literal, and then when it has two, and so forth, until we have enumerated all 1-simplified templates. This way, we are heuristically choosing the well-founded chains of smaller lengths while searching for a minimal template.

The second modification we introduce has the purpose of keeping  $C$  small in Algorithm 2. While computing a minimal

template for a single program, we first check whether any of the previously computed minimal templates is perhaps minimal for the new program. If so then we continue by analyzing the next program. As a result, the number of distinct minimal templates in *mintS* reduces, and so does the size of  $C$ . Further, due to the definition of  $\gamma_{\text{proc}, p}$ , our algorithm can produce the same annotated program for different templates. In that case, the result of running Houdini+Corral for the first template is same as the result of running them on the other template. We therefore use the concretized annotations to cache Corral's outcome and reuse the results when possible.

We set a timeout of 8 hours for the algorithm. If the execution reaches a timeout, we simply return the current value of  $C$ . In practice, for SDV, the template inference needs to be performed just once in a release cycle; hence, one can devote much more time for inference.

## 5. EVALUATION

We implemented our algorithms in a tool called PROOFMINIMIZATION<sup>5</sup>. It accepts a list of annotated programs written in Boogie [19] as input and computes the minimal repository template. The shared vocabulary is automatically defined to be the set of global variables common to all input programs. The implementation is a simple wrapper, of less than 1000 lines of C# code, around the implementation of Corral [18].

We ran our experiments on a cluster of 6 identical server-class machines. Each of the servers had Intel Xeon CPUs 1.8 GHz, 64 GB RAM and 16 logical processors. The total CPU time of our experiments exceeded well over a month. We relied on parallelism extensively to produce results in a reasonable amount of time.

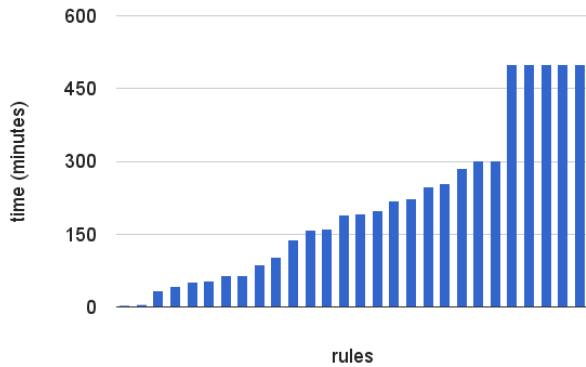
**Training Modes.** Our experiments evaluate two different ways of using the inferred templates. In the first mode, called MT, the inferred templates only augment the manual templates. This is achieved by adding the manual templates to each  $(p, t)$  pair in our program repository and requiring our algorithm to never throw out a manually-generated annotation. This mode of operation is more controlled: it does not seek to replace existing manual effort, but rather to just augment it.

The second mode of operation, called NT, does not use manual templates at all. This simulates the scenario when no manually-generated templates had been added to SDV. It answers how much of the manual effort behind the design of manual templates can be automated using our techniques.

Clearly, the quality of inferred templates will depend on the quality of the training set, i.e., the repository used for inferring a minimal repository template. Because it was never intended to use RTS for inferring templates, the training sets are sometimes inadequate for our approach. For instance, for a few rules, RTS only contains buggy drivers. (Inferring annotations from buggy programs is an interesting problem, but outside the scope of this paper.) Thus, we also evaluate expanding our training set by sampling a randomly-chosen fraction of the 66 real drivers and including them in the training set. The test set, i.e., the set of unseen drivers over which we evaluate the inferred templates then shrinks to the remaining drivers.

<sup>5</sup>Implementation is available at <https://github.com/boogie-org/corral/tree/master/AddOns/pminbench> along with supplemental material and examples that the reader can use to experiment.





**Figure 4: Annotation inference times for the 28 rules.**

Let  $\text{TRAIN}(f\%, m)$ , where  $0 \leq f \leq 100$  and  $m$  is either MT or NT, refer to the experiment where the training set consisted of all RTS drivers and  $f$  percent of the real drivers, and the inference was done in mode  $m$ . An exception is the special case of  $\text{TRAIN}(100\%, \text{MT})$  which denotes that all drivers were used in the training set as well as the test set. We use  $\text{TRAIN}(100\%, \text{MT})$  as a limit study on the quality of annotations that we can infer.

**Results.** For the experiment  $\text{TRAIN}(0\%, \text{NT})$ , the running time of `PROOFMINIMIZATION` is shown in Figure 4. It takes just a couple of minutes for some rules while it timeouts after 8 hours for five rules. Not including the rules for which `PROOFMINIMIZATION` times out, it takes roughly 2.5 hours on average to compute the result.

The results with various different training modes are reported in Tables 1 and 2. Each of the tables compares three versions of SDV. The version called “None” does not use any templates and captures the performance of Corral without any annotations. “Manual” refers to using manual templates, which is the currently-shipping production system. “Inferred” refers to using the set of annotations inferred by our tool. Each of the tables measures performance in terms of the number of timeouts (`#TO`), number of times the coverage bound was hit (`#Bnd`), the number of bugs reported (`#Bugs`), the average running time of Houdini+Corral (`Avg`), and the average running time of just Houdini (`Houd`).

Table 1(left) compares performance for  $\text{TRAIN}(0\%, \text{MT})$ . In this mode, our inferred set of annotations was empty for 12 of the 28 rules. The inferred set can be empty when the RTS wasn’t rich enough, or because the manual templates were sufficient. In this case, the performance of “Inferred” matches with that of “Manual”. Table 1 compares performance on the remaining 16 rules where we did infer annotations. These results demonstrate that our extra set of annotations are useful. The number of timeouts come down a fraction and the number of times the coverage bound was hit comes down significantly. All of these previously inconclusive cases, 54 in number, convert to a proved verdict. In summary, the total number of inconclusive answers drops down by 47%. Moreover, even though Houdini ran slower because of the extra annotations, the performance improvement in Corral made the overall system much faster (22%).

For  $\text{TRAIN}(30\%, \text{MT})$ , we did three runs, each time sampling a different fraction of the drivers. (The variance across

the three runs was very small, thus, we stopped with three runs.) The results for “Inferred” were averaged across the three runs and are shown in Table 1 (right). It is interesting that performance is similar to  $\text{TRAIN}(0\%, \text{MT})$ . Using a fraction of drivers did not provide much new information. However, the results of  $\text{TRAIN}(100\%, \text{MT})$  shown in Table 2 (left) do indicate that drivers potentially carry information not present in RTS. Learning over all drivers increased the quality of inferred annotations significantly (even though the running time of Houdini is highest in this setting).

Table 2 (right) shows results of  $\text{TRAIN}(0\%, \text{NT})$ . With no help from manual templates, Corral’s performance depends even more significantly on the inferred set of annotations. We were able to achieve a similar quality of results compared to using manual templates. There was a near-equal split between rules on which inferred annotations do better and ones on which manual templates do better. In consultation with the SDV team, we realized that the exercise of setting up manual templates often borrowed annotations *across* rules, where annotations useful for one rule would be generated for similar rules as well. In our setting, we did not explore sharing information between rules.

The number of abstract annotations per template, on average, was 12.78 for the manual templates. The  $\text{TRAIN}(0\%, \text{MT})$  experiment produced 13.63 annotations per template, on average. The  $\text{TRAIN}(100\%, \text{MT})$  experiment produced 16.6 annotations per template, on average. The  $\text{TRAIN}(0\%, \text{NT})$  experiment produced 3.3 annotations per template, on average, which is significantly smaller than the manual templates.

We supply the list of all rules and templates in the supplemental material. Here we summarize examples of useful annotations that our technique was able to infer, but were missed by experts. Firstly, for IRQL-based rules, the annotation `$floc == old($floc)` was necessary for establishing important loop invariants such as “the IRQL value is unchanged across loop iterations” (see example in Figure 3 and discussion in Sections 2.2 and 2.5). While the experts anticipated the latter (i.e., IRQL is unchanged) to be a useful annotation, they did not expect that an invariant over local variables would be required to prove it inductively.

A second instance of useful annotations found by our technique is for SDV rules with multiple model variables. Manual templates mostly have annotations that capture the effect of a procedure’s execution on a single model variable [28]. However, these are insufficient to capture inter-variable relationships. We infer several annotations over multiple variables which helped significantly for such rules.

To summarize,  $\text{TRAIN}(0\%, \text{MT})$  results show that we can significantly improve the performance of a production system by augmenting existing manual effort. The results for  $\text{TRAIN}(0\%, \text{NT})$  show that as new rules are developed and test cases are added to RTS, we can automatically generate useful annotations avoiding the need for further manual effort in coming up with new templates.

## 6. RELATED WORK

This work falls into the category of predicting program properties from codebases. For example, `JSNICE` learns from Javascript repositories on GitHub and predicts more legible identifier names and (unverified) type annotations [24]. In contrast, this work is the first attempt at inferring annotations from verification histories and demonstrating their use

**Table 1: Left: Results for Train(0%, MT) on a total of 16 rules with 873 verification instances. Right: Results for Train(30%, MT), averaged across three runs, on a total of 25 rules with 1002 verification instances.**

Config	#TO	#Bnd	#Bugs	Time (sec)	
				Avg	Houd
None	77	228	46	94.3	0
Manual	27	88	46	71.9	10.0
Inferred	24	37	46	51.6	12.1

Config	#TO	#Bnd	#Bugs	Time (sec)	
				Avg	Houd
None	112	265	64	104.3	0
Manual	50	91	64	70.5	10.1
Inferred	46.6	41	64	59.4	14.6

**Table 2: Left: Results for Train(100%, MT) on a total of 28 rules with 1420 verification instances. Right: Results for Train(0%, NT) on a total of 28 rules with 1420 verification instances.**

Config	#TO	#Bnd	#Bugs	Time (sec)	
				Avg	Houd
None	143	342	104	98.1	0
Manual	55	123	108	58.4	9.3
Inferred	45	25	108	46.1	16.2

Config	#TO	#Bnd	#Bugs	Time (sec)	
				Avg	Houd
None	143	342	104	98.1	0
Manual	55	123	108	58.44	9.3
Inferred	59	92	108	56.5	9.1

in an industrial-scale verification setting. Other approaches use codebases to predict different program properties rather than annotations [22, 23, 25]. For instance, work presented in [23] applies Bayesian optimization on existing codebases to learn a strategy for deciding for which part of an unseen program a static analyzer should sacrifice time for precision while performing the analysis.

There are verification and testing approaches that leverage previous versions of a program under analysis. For example, [11] improves performance of test generation for a program by leveraging existing tests belonging to a previous version. Regression verification verifies the equivalence of two successive versions of a program [12]. For similar programs, [12] argues that this verification task is easier than our goal here, i.e., formal verification of a stand alone program. A recent generalization of regression verification is differential assertion checking where a verifier checks that a bug is not introduced in going from one program version to another [14]. Techniques in [7] extrapolate from predicates used for verifying a program under sequential consistency to verify the same program under a relaxed memory model. Incremental verification attempts to reuse annotations corresponding to a program in the verification of its updated version [26, 8, 5]. In these works, the history typically consists of revisions of the program under analysis. We consider histories with programs that use the same kernel API but can be very different otherwise. For example, we can use annotations inferred from a storage and a modem driver to help verify an IEEE-1394 bus driver. Also, our technique uses the minimization mechanism to reduce the size of the accumulated annotations in order to retain practicality.

The techniques described in this paper can be used to infer invariants using verification histories. Previous approaches to invariant inference perform analysis only over the program under consideration. These include invariant inference using static analysis [6] or learning from concrete executions [10, 27]. Our work is complementary to these.

## 7. CONCLUSION

We present an algorithm for computing a *small* set of *useful* program annotations from a repository of past verification runs. We used this algorithm to improve the performance of SDV by generating better quality annotations than those produced by human experts. By utilizing the inferred an-

notations, we reduce the number of inconclusive answers by 47% while running 22% faster on average, even for a heavily-optimized system.

Our approach benefited from SDV’s separation of verifying driver correctness into checking multiple rules. By concentrating on one rule at a time, our inferred annotations were tailored to a specific property of drivers. In future work, we wish to study the generality of our algorithm in other verification domains that rely on annotations for constructing program proofs.

## 8. ACKNOWLEDGMENTS

The authors would like to thank Zilong Wang and Subhajit Roy for their help during initial development of the ideas presented in this paper, and Kenneth McMillan for his help with using the Duality verification engine. This work was in part supported by the National Science Foundation under grant CCF-1350574.

## 9. REFERENCES

- [1] A. Albarghouthi, A. Gurfinkel, Y. Li, S. Chaki, and M. Chechik. Ufo: Verification with interpolants and abstract interpretation. In *Tools and Algorithms for the Construction and Analysis of Systems*, pages 637–640. Springer, 2013.
- [2] T. Ball, V. Levin, and S. K. Rajamani. A decade of software model checking with SLAM. *Communications of the ACM*, 54(7):68–76, 2011.
- [3] T. Ball, R. Majumdar, T. D. Millstein, and S. K. Rajamani. Automatic predicate abstraction of C programs. In *Proceedings of the 2001 ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI), Snowbird, Utah, USA, June 20-22, 2001*, pages 203–213, 2001.
- [4] N. E. Beckman, A. V. Nori, S. K. Rajamani, R. J. Simmons, S. Tetali, and A. V. Thakur. Proofs from tests. *IEEE Trans. Software Eng.*, 36(4):495–508, 2010.
- [5] D. Beyer, S. Löwe, E. Novikov, A. Stahlbauer, and P. Wendler. Precision reuse for efficient regression verification. In *Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software*

- Engineering, ESEC/FSE'13, Saint Petersburg, Russian Federation, August 18-26, 2013, pages 389–399, 2013.
- [6] P. Cousot and R. Cousot. Abstract interpretation: A unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Conference Record of the Fourth ACM Symposium on Principles of Programming Languages, Los Angeles, California, USA, January 1977*, pages 238–252, 1977.
- [7] A. M. Dan, Y. Meshman, M. T. Vechev, and E. Yahav. Predicate abstraction for relaxed memory models. In *Static Analysis - 20th International Symposium, SAS 2013, Seattle, WA, USA, June 20-22, 2013. Proceedings*, pages 84–104, 2013.
- [8] G. Fedyukovich, A. Gurfinkel, and N. Sharygina. Incremental verification of compiler optimizations. In *NASA Formal Methods - 6th International Symposium, NFM 2014, Houston, TX, USA, April 29 - May 1, 2014. Proceedings*, pages 300–306, 2014.
- [9] C. Flanagan and K. R. M. Leino. Houdini, an annotation assistant for esc/java. In *FME 2001: Formal Methods for Increasing Software Productivity, International Symposium of Formal Methods Europe, Berlin, Germany, March 12-16, 2001, Proceedings*, pages 500–517, 2001.
- [10] P. Garg, C. Löding, P. Madhusudan, and D. Neider. ICE: A robust framework for learning invariants. In *Computer Aided Verification - 26th International Conference, CAV 2014, Held as Part of the Vienna Summer of Logic, VSL 2014, Vienna, Austria, July 18-22, 2014. Proceedings*, pages 69–87, 2014.
- [11] P. Godefroid, S. K. Lahiri, and C. Rubio-González. Statically validating must summaries for incremental compositional dynamic test generation. In *Static Analysis - 18th International Symposium, SAS 2011, Venice, Italy, September 14-16, 2011. Proceedings*, pages 112–128, 2011.
- [12] B. Godlin and O. Strichman. Regression verification: proving the equivalence of similar programs. *Softw. Test., Verif. Reliab.*, 23(3):241–258, 2013.
- [13] T. A. Henzinger, R. Jhala, R. Majumdar, and K. L. McMillan. Abstractions from proofs. In *Proceedings of the 31st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2004, Venice, Italy, January 14-16, 2004*, pages 232–244, 2004.
- [14] S. K. Lahiri, K. L. McMillan, R. Sharma, and C. Hawblitzel. Differential assertion checking. In *Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering, ESEC/FSE'13, Saint Petersburg, Russian Federation, August 18-26, 2013*, pages 345–355, 2013.
- [15] A. Lal and S. Qadeer. Powering the static driver verifier using corral. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering, (FSE-22), Hong Kong, China, November 16 - 22, 2014*, pages 202–212, 2014.
- [16] A. Lal and S. Qadeer. A program transformation for faster goal-directed search. In *Formal Methods in Computer-Aided Design, FMCAD 2014, Lausanne, Switzerland, October 21-24, 2014*, pages 147–154, 2014.
- [17] A. Lal and S. Qadeer. DAG inlining: a decision procedure for reachability-modulo-theories in hierarchical programs. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation, Portland, OR, USA, June 15-17, 2015*, pages 280–290, 2015.
- [18] A. Lal, S. Qadeer, and S. K. Lahiri. A solver for reachability modulo theories. In *Computer Aided Verification - 24th International Conference, CAV 2012, Berkeley, CA, USA, July 7-13, 2012 Proceedings*, pages 427–443, 2012.
- [19] K. R. M. Leino. This is boogie 2. *Manuscript KRML*, 178:131, 2008. <http://https://github.com/boogie-org/boogie>.
- [20] K. L. McMillan and A. Rybalchenko. Computing relational fixed points using interpolation. Technical report, Technical report, 2012. available from authors, 2013.
- [21] Microsoft. The Static Driver Verifier. [http://msdn.microsoft.com/en-us/library/windows/hardware/ff552808\(v=vs.85\).aspx](http://msdn.microsoft.com/en-us/library/windows/hardware/ff552808(v=vs.85).aspx).
- [22] A. Mishne, S. Shoham, and E. Yahav. Typestate-based semantic code search over partial programs. In *Proceedings of the 27th Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2012, part of SPLASH 2012, Tucson, AZ, USA, October 21-25, 2012*, pages 997–1016, 2012.
- [23] H. Oh, H. Yang, and K. Yi. Learning a strategy for adapting a program analysis via bayesian optimisation. In *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2015, part of SLASH 2015, Pittsburgh, PA, USA, October 25-30, 2015*, pages 572–588, 2015.
- [24] V. Raychev, M. T. Vechev, and A. Krause. Predicting program properties from big code. In *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2015, Mumbai, India, January 15-17, 2015*, pages 111–124, 2015.
- [25] V. Raychev, M. T. Vechev, and E. Yahav. Code completion with statistical language models. In *ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '14, Edinburgh, United Kingdom - June 09 - 11, 2014*, page 44, 2014.
- [26] O. Sery, G. Fedyukovich, and N. Sharygina. Incremental upgrade checking by means of interpolation-based function summaries. In *Formal Methods in Computer-Aided Design, FMCAD 2012, Cambridge, UK, October 22-25, 2012*, pages 114–121, 2012.
- [27] R. Sharma and A. Aiken. From invariant checking to invariant inference using randomized search. In *Computer Aided Verification - 26th International Conference, CAV 2014, Held as Part of the Vienna Summer of Logic, VSL 2014, Vienna, Austria, July 18-22, 2014. Proceedings*, pages 88–105, 2014.
- [28] R. S. Zvonimir Pavlinovic, Akash Lal. Inferring annotations for device drivers from verification histories. Technical report, Microsoft Research, April 2016.