

# Invariant Synthesis for Incomplete Verification Engines

Daniel Neider<sup>1</sup>, Pranav Garg<sup>2</sup>, P. Madhusudan<sup>3</sup>, Shambwaditya Saha<sup>3</sup>, and Daejun Park<sup>3</sup>

<sup>1</sup> Max Planck Institute for Software Systems, Kaiserslautern, Germany

<sup>2</sup> Amazon India, Bangalore, India

<sup>3</sup> University of Illinois at Urbana-Champaign, Champaign, IL, USA

**Abstract.** We propose a framework for synthesizing inductive invariants for incomplete verification engines, which soundly reduce logical problems in undecidable theories to decidable theories. Our framework is based on the counter-example guided inductive synthesis principle (CEGIS) and allows verification engines to communicate *non-provability information* to guide invariant synthesis. We show precisely how the verification engine can compute such non-provability information and how to build effective learning algorithms when invariants are expressed as Boolean combinations of a fixed set of predicates. Moreover, we evaluate our framework in two verification settings, one in which verification engines need to handle quantified formulas and one in which verification engines have to reason about heap properties expressed in an expressive but undecidable separation logic. Our experiments show that our invariant synthesis framework based on non-provability information can both effectively synthesize inductive invariants and adequately strengthen contracts across a large suite of programs.

## 1 Introduction

The paradigm of *deductive verification* [22,31] combines manual annotations and semi-automated theorem proving to prove programs correct. Programmers annotate code they develop with contracts and inductive invariants, and use high-level directives to an underlying, mostly-automated logic engine to verify their programs correct. Several mature tools have emerged that support such verification, in particular tools based on the intermediate verification language BOOGIE [3] and the SMT solver Z3 [45] (e.g., VCC [13] and DAFNY [40]). Various applications that use such tools to prove systems correct using manual annotations have been developed, including Microsoft Hypervisor verification [14], reliable systems code such as VERVE [56], ExpressOS [42], and Ironclad apps [30], as well as distributed systems in IronFleet [29]. Fully automated use of such engines for shallow specifications have also emerged, such as CORRAL [38] for verifying device drivers, CST [10] to certify transactions in online services, and GPUVERIFY [4] to ensure race-freedom in GPU kernels.

Viewed through the lens of deductive verification, the primary challenges in automating verification are two-fold. First, even when strong annotations in terms of contracts and inductive invariants are given, the validity problem for the resulting verification conditions is often undecidable (e.g., in reasoning about the heap, reasoning with quantified logics, and reasoning with non-linear arithmetic). Second, the synthesis of loop invariants and strengthenings of contracts that prove a program correct needs to be automated so as to lift this burden currently borne by the programmer.

A standard technique to solve the first problem (i.e., intractability of validity checking of verification conditions) is to build automated, sound-but-incomplete verification engines for validating verification conditions, thus skirting the undecidability barrier. Several such techniques exist; for instance, for reasoning with quantified formulas, tactics such as E-matching [17,44], pattern-based quantifier instantiation [17], and model-based quantifier instantiation [26] are effective in practice, though they are not complete for most background theories. In the realm of heap verification, the so-called *natural proof method* explicitly aims to provide automated and sound-but-incomplete methods for checking validity of verification conditions with specifications in separation logic [50,48,12]. This method searches for proofs based on induction on recursively defined data structures, which is reduced to validity problems in *decidable logics with quantification* that enables an efficient search for such proofs using SMT solvers.

Turning to the second problem of invariant generation, several techniques have emerged that can synthesize invariants automatically when validation of verification conditions fall in decidable classes. Prominent among these are interpolation [43] and IC3/PDR [6,19]. These techniques generalize from information gathered in proving underapproximations of the program correct and are quite effective [5]—their efficacy in dealing with programs where the underlying logics are undecidable, however, is unclear. Moreover, a class of counterexample guided inductive synthesis (CEGIS) methods have emerged recently, including the ICE learning model [24] for which various instantiations exist [24,52,25,37]. The key feature of the latter methods is a program-agnostic, data-driven learner that learns invariants in tandem with a verification engine that provides concrete program configurations as counterexamples to incorrect invariants.

Although classical invariant synthesis techniques, such as HOUDINI [21], are sometimes used with incomplete verification engines, to the best of our knowledge there is no fundamental argument as to why this should work in general. In fact, we are not aware of any systematic technique for synthesizing invariants when the underlying verification problem falls in an undecidable theory. When verification is undecidable and the engine resorts to sound but incomplete heuristics to check validity of verification conditions, it is unclear how to extend interpolation/IC3/PDR techniques to this setting. Data-driven learning of invariants is also hard to extend since the verification engine typically cannot generate a concrete model for the negation of verification conditions when verification fails. Hence, it cannot produce the concrete configurations that the learner needs.

**The main contribution of this paper is a general, learning-based invariant synthesis framework that learns invariants using non-provability information provided by verification engines.** Intuitively, when a conjectured invariant results in verification conditions that cannot be proven, the idea is that the verification engine must return information that generalizes the reason for non-provability, hence pruning the space of future conjectured invariants. Our framework assumes a verification engine for an undecidable theory  $\mathcal{U}$  that reduces verification conditions to a decidable theory  $\mathcal{D}$  (e.g., using heuristics such as bounded quantifier instantiation to remove universal quantifiers, function unfolding to remove recursive definitions, and so on) that permits producing models for satisfiable formulas. The translation is assumed to be conservative in the sense that if the translated formula in  $\mathcal{D}$  is valid, then we are assured that the original verification condition is  $\mathcal{U}$ -valid. If the verification condition is found to be not  $\mathcal{D}$ -valid (i.e., its negation is satisfiable), on the other hand, our framework describes how to extract non-provability information from the  $\mathcal{D}$ -model. This information is encoded as conjunctions and disjunctions in a Boolean theory  $\mathcal{B}$ , called *conjunctive/disjunctive non-provability information (CD-NPI)*, and communicated back to the learner. To complete our framework, we show how the formula-driven problem of learning expressions from CD-NPI constraints can be reduced to the data-driven ICE model. This reduction allows us to use a host of existing ICE learning algorithms and results in a robust invariant synthesis framework that guarantees to synthesize a provable invariant if one exists. We present the framework in Section 2 in detail.

However, our CD-NPI learning framework has non-trivial requirements on the verification engine, and building (or adapting) appropriate engines is not straightforward. To show that our framework is indeed applicable and effective in practice, **our second contribution is the application of our technique to two real-world verification settings.**

The first setting, presented in Section 3, is the verification of dynamically manipulated data-structures against rich logics that combine properties of structure, separation, arithmetic, and data—an important problem where verification often falls in undecidable theories. We show how *natural proof verification engines* [48], which are essentially sound-but-incomplete verification engines that translate a powerful undecidable separation logic called DRYAD to decidable logics, can be fit into our framework. We then implement a prototype of such a natural proof verification engines on top of the program verifier BOOGIE [3] and demonstrate that this prototype is able to fully automatically verify a large suite of benchmarks, containing standard algorithms for manipulating singly and doubly linked lists, sorted lists, as well as balanced and sorted trees. Automatically synthesizing invariants for this suite of heap manipulating programs against an expressive separation logic is very challenging, and we do not know of any other current technique that can automatically prove all of them. Thus, we have to leave a comparison to other approaches for future work.

The second setting is the verification of programs against specifications with universal quantification, which occur, for instance, when defining recursive properties. Again, we implement a prototype over BOOGIE and

demonstrate its effectiveness on a series of benchmarks taken from verification competitions and real-world systems. We describe this application in Section 4 and conclude in Section 5.

## Related Work

Techniques for invariant synthesis include abstract interpretation [16], interpolation [43], IC3 [6], predicate abstraction [2], abductive inference [18], as well as synthesis algorithms that rely on constraint solving [27,28,15]. Complementing them are data-driven invariant synthesis techniques based on learning, such as Daikon [20] that learn likely invariants, and HOUDINI [21] and ICE [24] that learn inductive invariants. The latter typically requires a teacher that can generate counter-examples if the conjectured invariant is not adequate or inductive. Classically, this is possible only when the verification conditions of the program fall in decidable logics. In this paper, we investigate data-driven invariant synthesis for incomplete verification engines and show that the problem can be reduced to ICE learning if the learning algorithm learns from non-provability information and produces hypotheses in a class that is restricted to positive Boolean formulas over a fixed set of predicates. Data-driven synthesis of invariants has regained recent interest [55,53,54,23,24,52,37,57,47,46] and our work addresses an important problem of synthesizing invariants for programs whose verifications conditions fall in undecidable fragments.

Our application to learning invariants for heap programs builds upon DRYAD [50,48], and the natural proof technique line of work for heap verification developed by Qiu et al. Techniques, similar to DRYAD, for automated reasoning of dynamically manipulated data structure programs have also been proposed in [12,11]. However, unlike our current work, none of these works synthesize heap invariants. Given invariant annotations in their respective logics, they provide procedures to validate if the verification conditions are valid. There has also been a lot of work on synthesizing invariants for separation logic using shape analysis [51,9,39]. However, most of them are tailored for memory safety and shallow properties rather than rich properties that check full functional correctness of data structures. Interpolation has also been suggested recently to synthesize invariants involving a combination of data and shape properties [1]. It is, however, not clear how the technique can be applied to a more complicated heap structure, such as an AVL tree, where shape and data properties are not cleanly separated but are intricately connected. Recent work also includes synthesizing heap invariants in the logic from [32] by extending IC3 [33,34].

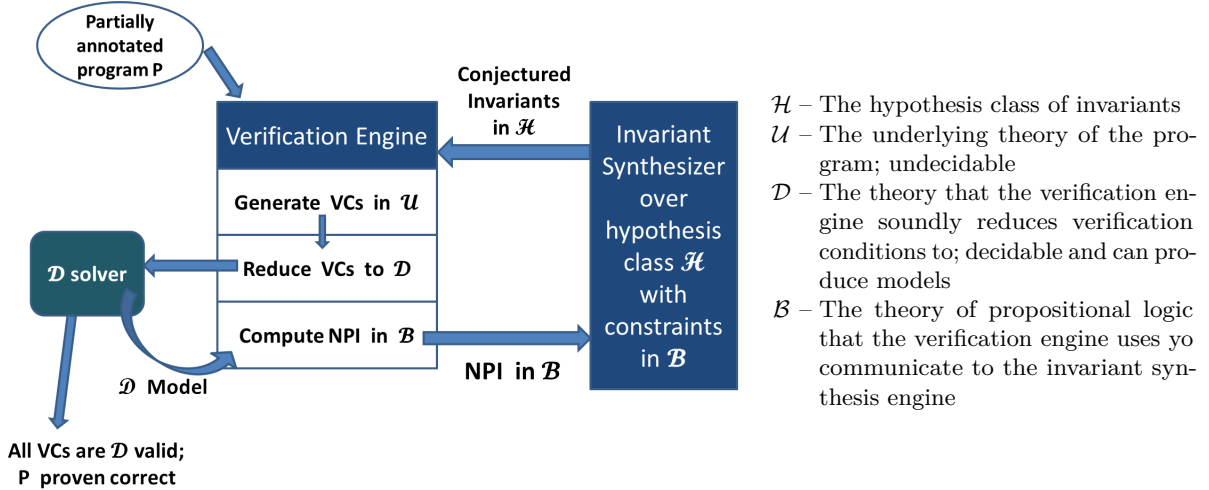
In this work, our learning algorithm synthesizes invariants over a fixed set of predicates. When all programs belong to a specific class, such as the class of programs manipulating data structures, these predicates can be uniformly chosen using templates. Investigating automated ways for discovering candidate predicates is a very interesting future direction. Related work in this direction includes recent works [47,46].

## 2 An Invariant Synthesis Framework for Incomplete Verification Engines

In this section, we develop our framework for synthesizing inductive invariants for incomplete verification engines, using a counter-example guided inductive synthesis approach. We do this in the setting where the hypothesis space consists of formulas that are Boolean combinations of a fixed set of predicates  $\mathcal{P}$ , which need not be finite for the general framework—when developing concrete learning algorithms later, we will assume  $\mathcal{P}$  is a finite set of predicates. For the rest of this section, let us fix a program  $P$  that is annotated with assertions (and possibly with some partial annotations describing pre-conditions, post-conditions, and assertions). Moreover, we refer to a formula  $\alpha$  being weaker (or stronger) than  $\beta$  in a logic  $\mathcal{L}$ , and by this we mean that  $\vdash_{\mathcal{L}} \beta \Rightarrow \alpha$  (or  $\vdash_{\mathcal{L}} \alpha \Rightarrow \beta$ ), respectively, where  $\vdash_{\mathcal{L}} \varphi$  means that  $\varphi$  is valid in  $\mathcal{L}$ .

Figure 1 (on Page 4) depicts our general framework of invariant synthesis when verification is undecidable. We fix several parameters for our verification effort. First, let us assume a uniform signature for logic, in terms of constant symbols, relation symbols, functions, and types. We will, for simplicity of exposition, use the same syntactic logic for the various logics  $\mathcal{U}$ ,  $\mathcal{D}$ ,  $\mathcal{B}$  in our framework as well as for the logic  $\mathcal{H}$  used to express invariants.

Let us fix  $\mathcal{U}$  as the underlying theory that is ideally needed for validating the verification conditions that arise for the program; we presume validity of formulas in  $\mathcal{U}$  is undecidable. Since  $\mathcal{U}$  is an undecidable



**Fig. 1.** A non-provability information (NPI) framework for invariant synthesis when the verification logic is undecidable

theory, the engine will resort to sound approximations (e.g., using bounded quantifier instantiations using mechanisms such as triggers [44], bounded unfolding of recursive functions, or natural proofs [48]) to reduce this logical task to a *decidable* theory  $\mathcal{D}$ . This reduction is assumed to be sound in the sense that if the resulting formulas in  $\mathcal{D}$  are valid, then the verification conditions are valid in  $\mathcal{U}$  as well. If a formula is found *not valid* in  $\mathcal{D}$ , then we require that the logic solver for  $\mathcal{D}$  returns a model for the negation of the formula.<sup>4</sup> Note that this model may not be a model for the negation of the formula in  $\mathcal{U}$ .

Moreover, we fix a hypothesis class  $\mathcal{H}$  for invariants consisting of *positive* Boolean combination of predicates in a fixed set of predicates  $\mathcal{P}$ . Note that restricting to *positive* formulas over  $\mathcal{P}$  is not a restriction, as one can always add negations of predicates to  $\mathcal{P}$ , thus effectively synthesizing any Boolean combination of predicates. The restriction to positive Boolean formulas is in fact desirable, as it allows restricting invariants to *not* negate certain predicates, which is useful when predicates have intuitionistic definitions (as several recursive definitions of heap properties do).

The invariant synthesis proceeds in rounds, where in each round the synthesizer proposes invariants in  $\mathcal{H}$ . The verification engine generates verification conditions in accordance to these invariants in the underlying theory  $\mathcal{U}$ . It then proceeds to translate them into the decidable theory  $\mathcal{D}$ , and gives them to a solver that decides validity in the theory  $\mathcal{D}$ . If the verification conditions are found to be  $\mathcal{D}$ -valid, then by virtue of the fact that the verification engine reduced VCs in a sound fashion to  $\mathcal{D}$ , we are done proving the program  $P$ .

However, if the formula is found not to be  $\mathcal{D}$ -valid, the solver returns a  $\mathcal{D}$ -model for the negation of the formula. The verification engine then extracts from this model certain *non-provability information (NPI)*, expressed as Boolean formulas in a Boolean theory  $\mathcal{B}$ , that captures more general reasons why the verification failed and eliminates not only the current conjectured invariant but also others that can be inferred to be incorrect from the current verification effort (the rest of this section is devoted to developing this notion of non-provability information). This non-provability information is communicated to the synthesizer, which then proceeds to synthesize a new conjecture invariant that satisfies the non-provability constraints provided in all previous rounds. The following example illustrates the logics involved in our framework in the context heap-manipulating programs.

*Example 1.* In a verification setting involving heaps, the logic  $\mathcal{U}$  could be a rich separation logic with recursive definitions and  $\mathcal{D}$  could be the quantifier-free theory of uninterpreted functions, arithmetic, and sets. The

<sup>4</sup> Note that our framework requires model construction in the theory  $\mathcal{D}$ . Hence, incomplete logic solvers for  $\mathcal{U}$  that simply time out after some time threshold or search for a proof of a particular kind and give up otherwise are not suitable candidates.

verification engine can reduce verification conditions in  $\mathcal{U}$  to  $\mathcal{D}$  by partially unfolding recursive definitions and expressing heaplets using sets, to obtain sound but incomplete automatic validity checking.

Futhermore, the theory  $\mathcal{B}$  can be chosen to be the just the propositional theory over a set of predicates  $\mathcal{P}$ . The verification engine will then communicate formulas over  $\mathcal{B}$  to the synthesis engine that restrict the class of invariants such that the synthesis engine can generate in the future.  $\triangleleft$

In order for the verification engine to extract meaningful non-provability information, we make the following natural assumption, called *normality*, which essentially states that the engine can do at least some minimal Boolean reasoning.

**Definition 1.** *A verification engine is normal if it satisfies two properties:*

1. *if the engine cannot prove the validity of the Hoare triple  $\{\alpha\}s\{\gamma\}$  and  $\vdash_{\mathcal{B}} \delta \Rightarrow \gamma$ , then it cannot prove the validity of the Hoare triple  $\{\alpha\}s\{\delta\}$ ; and*
2. *if the engine cannot prove the validity of the Hoare triple  $\{\gamma\}s\{\beta\}$  and  $\vdash_{\mathcal{B}} \gamma \Rightarrow \delta$ , then it cannot prove the validity of the Hoare triple  $\{\delta\}s\{\beta\}$ .*

Intuitively, Condition 1 of Definition 1 means that if an oracle cannot prove the validity of  $\{\alpha\}s\{\gamma\}$ , then it cannot prove the validity of any strengthening  $\delta$  of the postcondition  $\gamma$ . Similarly, Condition 2 means that if an oracle cannot prove the validity of  $\{\gamma\}s\{\beta\}$ , then it cannot prove the validity of any weakening  $\delta$  of the precondition  $\gamma$ .

The remainder of this section is now structured as follows. In Section 2.1, we first develop an appropriate language to communicate non-provability constraints, which allow the learner to appropriately weaken or strengthen a future hypothesis. It turns out that *pure conjunctions* and *pure disjunctions* over  $\mathcal{P}$ , which we term *CD-NPI constraints* (conjunctive/disjunctive non-provability information constraints), are sufficient for this purpose. We also describe concretely how the verification engine can extract this non-provability information from  $\mathcal{D}$ -models that witness that negations of VCs are satisfiable. Then, in Section 2.2, we show how to build learners for CD-NPI constraints by reducing this learning problem to another, well-studied learning framework for invariants called ICE learning. We illustrate our framework with an example in Section 2.3 and finally argue in Section 2.4 that our framework is sound and guarantees to converge to a provable invariant if one exists.

## 2.1 Conjunctive/Disjunctive Non-provability Information

We assume that the underlying decidable theory  $\mathcal{D}$  is stronger than propositional theory  $\mathcal{B}$ , meaning that every valid statement in  $\mathcal{B}$  is valid in  $\mathcal{D}$  as well. The reader may want to keep the following as a running example where  $\mathcal{D}$  is the decidable theory of uninterpreted functions and linear arithmetic, say. In this setting, a formula is  $\mathcal{B}$ -valid if, when treating atomic formulas as Boolean variables, the formula is propositionally valid. For instance,  $f(x) = y \Rightarrow f(f(x)) = f(y)$  will not be  $\mathcal{B}$ -valid though it is  $\mathcal{D}$ -valid, while  $f(x) = y \vee \neg(f(x) = y)$  is  $\mathcal{B}$ -valid.

To formally define CD-NPI constraints and their extraction from a failed verification attempt, let us first introduce the following notation. For any  $\mathcal{U}$ -formula  $\varphi$ , let  $approx(\varphi)$  denote the  $\mathcal{D}$ -formula that the verification engine generates such that the  $\mathcal{D}$ -validity of  $approx(\varphi)$  implies the  $\mathcal{U}$ -validity of  $\varphi$ . Moreover, for any Hoare triple  $\{\alpha\}s\{\beta\}$ , let  $VC(\{\alpha\}s\{\beta\})$  denote the verification condition corresponding to the Hoare triple that the verification engine generates.

Let us now assume, for the sake of a simpler exposition, that the program has a single annotation hole  $A$  where we need to synthesize an inductive invariant and prove the program correct. Further, suppose the learner conjectures an annotation  $\gamma$  as an inductive invariant for the annotation hole  $A$ , and the verification engine fails to prove the verification condition corresponding to a Hoare triple  $\{\alpha\}s\{\beta\}$ , where either  $\alpha$ ,  $\beta$ , or both could involve the synthesized annotation. This means that the negation of  $approx(VC(\{\alpha\}s\{\beta\}))$  is  $\mathcal{D}$ -satisfiable and the verification engine needs to extract non-provability information from a model of it. To this end, we assume that every program snippet  $s$  has been augmented with a set of ghost variables  $g_1, \dots, g_n$  that track the predicates  $p_1, \dots, p_n$  mentioned in the invariant (i.e., these ghost variables are

assigned the values of the predicates). The valuation  $\mathbf{v} = \langle v_1, \dots, v_n \rangle$  of the ghost variables in the model before the execution of  $s$  and the valuation  $\mathbf{v}' = \langle v'_1, \dots, v'_n \rangle$  after the execution of  $s$  can then be used to derive non-provability information, as we describe shortly.

The type of non-provability information the verification engine extracts depends on where the annotation appears in a Hoare triple  $\{\alpha\}s\{\beta\}$ . More specifically, the synthesized annotation might appear in  $\alpha$ , in  $\beta$ , or in both. We now handle all three cases individually.

- Assume the verification of a Hoare triple of the form  $\{\alpha\}s\{\gamma\}$  fails (i.e., the verification engine cannot prove a verification condition where the pre-condition  $\alpha$  is a user-supplied annotation and the post-condition is the synthesized annotation  $\gamma$ ). Then,  $\text{approx}(VC(\{\alpha\}s\{\gamma\}))$  is not  $\mathcal{D}$ -valid, and the decision procedure for  $\mathcal{D}$  would generate a model for its negation.

Since  $\gamma$  is a positive Boolean combination, the reason why  $\mathbf{v}'$  does not satisfy  $\gamma$  is due to the variables mapped to *false* by  $\mathbf{v}'$ , as any valuation extending this will not satisfy  $\gamma$ . Intuitively, this means that the  $\mathcal{D}$ -solver is not able to prove the predicates in  $P_{false} = \{p_i \mid v'_i = false\}$ . In other words,  $\{\alpha\}s\{\bigvee P_{false}\}$  is unprovable (a witness to this fact is the model of the negation of  $\text{approx}(VC(\{\alpha\}s\{\gamma\}))$  from which the values  $\mathbf{v}'$  are derived). Note that any invariant  $\gamma'$  that is stronger than  $\bigvee P_{false}$  will result in an unprovable VC due to the verification engine being normal. Consequently we can choose  $\chi = \bigvee P_{false}$  as the weakening constraint, demanding that future invariants should not be stronger than  $\chi$ .

The verification engine now communicates  $\chi$  to the synthesizer, asking it never to conjecture in future rounds invariants  $\gamma''$  that are stronger than  $\chi$  (i.e., such that  $\not\vdash_{\mathcal{B}} \gamma'' \Rightarrow \chi$ ).

- The next case is when a Hoare triple of the form  $\{\gamma\}s\{\beta\}$  fails to be proven (i.e., the verification engine cannot prove a verification condition where the post-condition  $\beta$  is a user-supplied annotation and the pre-condition is the synthesized annotation  $\gamma$ ). Using similar arguments as above, the *conjunction*  $\eta = \bigwedge \{p_i \mid v_i = true\}$  of the predicates mapped to *true* by  $\mathbf{v}$  in the corresponding  $\mathcal{D}$ -model gives a stronger precondition  $\eta$  such that  $\{\eta\}s\{\alpha\}$  is not provable. Hence,  $\eta$  is a valid *strengthening* constraint. The verification engine now communicates  $\eta$  to the synthesizer, asking it never to conjecture in future rounds invariants  $\gamma''$  that are weaker than  $\eta$  (i.e., such that  $\not\vdash_{\mathcal{B}} \eta \Rightarrow \gamma''$ ).
- Finally, consider the case when the Hoare triple is of the form  $\{\gamma\}s\{\gamma\}$  and fails to be proven (i.e., the verification engine cannot prove a verification condition where the pre- and post-condition is the synthesized annotation  $\gamma$ ). In this case, the verification engine can offer advice on how  $\gamma$  can be strengthened or weakened to avoid this model. Analogous to the two cases above, the verification engine extracts a pair of formulas  $(\eta, \chi)$ , called an *inductivity constraint*, based on the variables mapped to *true* by  $\mathbf{v}$  and to *false* by  $\mathbf{v}'$ . The meaning of such a constraint is that the invariant synthesizer must conjecture in future rounds invariants  $\gamma''$  such that either  $\not\vdash_{\mathcal{B}} \eta \Rightarrow \gamma''$  or  $\not\vdash_{\mathcal{B}} \gamma'' \Rightarrow \chi$  holds.

This leads to the following scheme, where  $\gamma$  denotes the conjectured invariant:

- When a Hoare triple of the form  $\{\alpha\}s\{\gamma\}$  fails, the verification engine returns the  $\mathcal{B}$ -formula

$$\bigvee_{i \mid v'_i = false} p_i$$

as a *weakening constraint*.

- When a Hoare triple of the form  $\{\gamma\}s\{\beta\}$  fails, the verification engine returns the  $\mathcal{B}$ -formula

$$\bigwedge_{i \mid v_i = true} p_i$$

as a *strengthening constraint*.

- When a Hoare triple of the form  $\{\gamma\}s\{\gamma\}$  fails, the verification engine returns the pair

$$\left( \bigwedge_{i \mid v_i = true} p_i, \bigvee_{i \mid v'_i = false} p_i \right)$$

of  $\mathcal{B}$ -formulas as an inductivity constraint.

It is not hard to verify that the above formulas are proper strengthening and weakening constraints, in the sense that *any* inductive invariant must satisfy these constraints. This motivates the following form of non-provability information.

**Definition 2 (CD-NPI Samples).** Let  $\mathcal{P}$  be a set of predicates. A CD-NPI sample (short for conjunction-disjunction-NPI sample) is a triple  $\mathfrak{S} = (W, S, I)$  consisting of

- a finite set  $W$  of disjunctions over  $\mathcal{P}$  (weakening constraints);
- a finite set  $S$  of conjunctions over  $\mathcal{P}$  (strengthening constraints); and
- a finite set  $I$  of pairs, where the first element is a conjunction and the second is a disjunction over  $\mathcal{P}$  (inductivity constraints).

An annotation  $\gamma$  is consistent with a CD-NPI sample  $\mathfrak{S} = (W, S, I)$  if

- $\not\vdash_{\mathcal{B}} \gamma \Rightarrow \chi$  for each  $\chi \in W$ ;
- $\not\vdash_{\mathcal{B}} \eta \Rightarrow \gamma$  for each  $\eta \in S$ ; and
- $\not\vdash_{\mathcal{B}} \eta \Rightarrow \gamma$  or  $\not\vdash_{\mathcal{B}} \gamma \Rightarrow \chi$  for each  $(\eta, \chi) \in I$ .

A CD-NPI learner is an effective procedure that synthesizes, given an CD-NPI sample, an annotation  $\gamma$  consistent with the sample. In our framework, the process of proposing candidate annotations and checking them repeats until the learner proposes a valid annotation or it detects that no valid annotation exists (e.g., if the class of candidate annotations is finite and all annotations are exhausted). We comment on using an CD-NPI learner in this iterative fashion below.

## 2.2 Building CD-NPI Learners

Let us now turn to the problem of building efficient learning algorithms for CD-NPI constraints. To this end, we assume that the set of predicates  $\mathcal{P}$  is finite.

Roughly speaking, the CD-NPI learning problem is to synthesize annotations that are positive Boolean combinations of predicates in  $\mathcal{P}$  and that are consistent with given CD-NPI samples. Though this is a learning problem where samples are *formulas*, in this section we will reduce CD-NPI learning to a learning problem from *data*. In particular, we will show that CD-NPI learning reduces to the ICE learning framework for learning positive Boolean formulas. The latter is a well-studied framework, and the reduction allows us to use efficient learning algorithms developed for ICE learning in order to build CD-NPI learners.

We now first recap the ICE-learning framework and then reduce CD-NPI learning to ICE learning. Finally, we briefly sketch how the popular HOUDINI algorithm can be seen as an ICE learning algorithm, which, in turn, allows us to HOUDINI as an CD-NPI learning algorithm.

**The ICE learning framework** Although the ICE learning framework [24] is a general framework for learning inductive invariants, we consider here the case of learning Boolean formulas. To this end, let us fix a set  $B$  of Boolean variables, and let  $\mathcal{H}$  be a subclass of positive Boolean formulas over  $B$ , called the hypothesis class, which specifies the admissible solutions to the learning task.

The objective of the (passive) ICE learning algorithm is to learn a formula in  $\mathcal{H}$  from a sample of positive examples, negative examples, and implication examples. More formally, if  $\mathcal{V}$  is the set of valuations  $v: B \rightarrow \{\text{true}, \text{false}\}$  (mapping variables in  $B$  to true or false), then an *ICE sample* is a triple  $\mathcal{S} = (S_+, S_-, S_{\Rightarrow})$  where

- $S_+ \subseteq \mathcal{V}$  is a set of *positive examples*;
- $S_- \subseteq \mathcal{V}$  is a set of *negative examples*; and
- $S_{\Rightarrow} \subseteq \mathcal{V} \times \mathcal{V}$  is a set of *implications*.

Note that positive and negative examples are *concrete* valuations of the variables  $B$ , and the implication examples are pairs of such concrete valuations.

A formula  $\varphi$  is said to be *consistent with an ICE sample*  $\mathcal{S}$  if it satisfies the following three conditions:<sup>5</sup>

<sup>5</sup> In the following,  $\models$  denotes the usual satisfaction relation.

- $v \models \varphi$  for each  $v \in S_+$ ;
- $v \not\models \varphi$  for each  $v \in S_-$ ; and
- $v_1 \models \varphi$  implies  $v_2 \models \varphi$  for each  $(v_1, v_2) \in S_{\Rightarrow}$ .

In algorithmic learning theory, one distinguished between *passive learning* and *iterative learning*. The former refers to a learning setting in which a learning algorithm is confronted with a finite set of data and has to learn a concept that is consistent with this data. Using our terminology, the *passive ICE learning problem* for a hypothesis class  $\mathcal{H}$  is then

“given an ICE sample  $\mathcal{S}$ , find a formula in  $\mathcal{H}$  that is consistent with  $\mathcal{S}$ ”.

Recall that we here require the learner to learn positive Boolean formulas, which is slightly stricter than the original definition [24].

Iterative learning, on the other hand, is the iteration of passive learning where new data is added to the sample from one iteration to the next. In a verification context, this new data is generated by the verification engine in response to incorrect annotations and used to guide the learning algorithm towards an annotation that is adequate to prove the program. To reduce our learning framework to ICE learning, it is therefore sufficient to reduce the (passive) CD-NPI learning problem described above to the passive ICE learning problem. We do this next.

**Reduction of passive CD-NPI learning to passive ICE learning** Let  $\mathcal{H}$  be a subclass of positive Boolean formulas. We reduce the CD-NPI learning problem for  $\mathcal{H}$  to the ICE learning problem for  $\mathcal{H}$ . The main idea is to (a) treat each predicate  $p \in \mathcal{P}$  as a Boolean variable for the purpose of ICE learning and (b) to translate a CD-NPI sample  $\mathfrak{G}$  into an *equi-consistent* ICE sample  $\mathcal{S}_{\mathfrak{G}}$ , meaning that a positive Boolean formula is consistent with  $\mathfrak{G}$  if and only if it is consistent with  $\mathcal{S}_{\mathfrak{G}}$ . Then, learning a consistent formula in the CD-NPI framework for the hypothesis class  $\mathcal{H}$  reduces to learning consistent formulas in  $\mathcal{H}$  in the ICE learning framework.

The following lemma will help translate between the two frameworks. Its proof is straightforward, and follows from the fact that for any *positive* formula  $\alpha$ , if a valuation  $v$  sets a larger subset of propositions to true than  $v'$  does and  $v' \models \alpha$ , then  $v \models \alpha$  as well.

**Lemma 1.** *Let  $v$  be a valuation of  $\mathcal{P}$  and  $\alpha$  be a positive Boolean formula over  $\mathcal{P}$ . Then, the following holds:*

- $v \models \alpha$  if and only if  $\vdash_{\mathcal{B}} (\bigwedge_{p|v(p)=\text{true}} p) \Rightarrow \alpha$  (and, thus,  $v \not\models \alpha$  if and only if  $\not\vdash_{\mathcal{B}} (\bigwedge_{p|v(p)=\text{true}} p) \Rightarrow \alpha$ ).
- $v \models \alpha$  if and only if  $\not\vdash_{\mathcal{B}} \alpha \Rightarrow (\bigvee_{p|v(p)=\text{false}} p)$ .

This motivates our translation, which relies on two functions,  $d$  and  $c$ . The function  $d$  translates a disjunction  $\bigvee J$ , where  $J \subseteq \mathcal{P}$  is a subset of propositions, into the valuation

$$d(\bigvee J) = v \text{ with } v(p) = \text{false} \text{ if and only if } p \in J.$$

The function  $c$  translates a conjunction  $\bigwedge J$ , where  $J \subseteq \mathcal{P}$ , into the valuation

$$c(\bigwedge J) = v \text{ with } v(p) = \text{true} \text{ if and only if } p \in J.$$

By substituting  $v$  in Lemma 1 with  $c(\bigwedge J)$  and  $d(\bigvee J)$ , respectively, one immediately obtains the following result.

**Lemma 2.** *Let  $J \subseteq \mathcal{P}$  and  $\alpha$  be a positive Boolean formula over  $\mathcal{P}$ . Then, the following holds: (a)  $c(\bigwedge J) \not\models \alpha$  if and only if  $\not\vdash_{\mathcal{B}} \bigwedge J \Rightarrow \alpha$ , and (b)  $d(\bigvee J) \models \alpha$  if and only if  $\not\vdash_{\mathcal{B}} \alpha \Rightarrow \bigvee J$ .*

Based on the functions  $c$  and  $d$ , the translation of a CD-NPI sample into an equi-consistent ICE sample is as follows.



**Definition 3.** Given a CD-NPI sample  $\mathfrak{S} = (W, S, I)$ , the ICE sample  $\mathcal{S}_{\mathfrak{S}} = (S_+, S_-, S_{\Rightarrow})$  is defined by  $S_+ = \{d(\bigvee J) \mid \bigvee J \in W\}$ ,  $S_- = \{c(\bigwedge J) \mid \bigwedge J \in S\}$ , and  $S_{\Rightarrow} = \{(c(\bigwedge J_1), d(\bigvee J_2)) \mid (\bigwedge J_1, \bigvee J_2) \in I\}$ .

By virtue of the lemma above, we can now establish the correctness of the reduction from the CD-NPI learning problem to the ICE learning problem.

**Theorem 1.** Let  $\mathfrak{S} = (W, S, I)$  be a CD-NPI sample,  $\mathcal{S}_{\mathfrak{S}} = (S_+, S_-, S_{\Rightarrow})$  the ICE sample as in Definition 3,  $\gamma$  a positive Boolean formula over  $\mathcal{P}$ . Then,  $\gamma$  is consistent with  $\mathfrak{S}$  if and only if  $\gamma$  is consistent with  $\mathcal{S}_{\mathfrak{S}}$ .

*Proof.* Let  $\mathfrak{S} = (W, S, I)$  be an CD-NPI sample, and let  $\mathcal{S}_{\mathfrak{S}} = (S_+, S_-, S_{\Rightarrow})$  the ICE sample as in Definition 3. Moreover, let  $\gamma$  be a positive Boolean formula. We prove Theorem 1 by considering each weakening, strengthening, and inductivity constraint together with their corresponding positive, negative, and implication examples individually.

- Pick a weakening constraint  $\bigvee J \in W$ , and let  $v \in S_+$  with  $v = d(\bigvee J)$  be the corresponding positive sample. Moreover, assume that  $\gamma$  is consistent with  $\mathfrak{S}$  and, thus,  $\not\vdash_{\mathcal{B}} \gamma \Rightarrow \bigvee J$ . By Lemma 2, this is true if and only if  $d(\bigvee J) \models \gamma$ . Hence,  $v \models \gamma$ .  
Conversely, assume that  $\gamma$  is consistent with  $\mathcal{S}$ . Thus,  $v \models \gamma$ , which means  $d(\bigvee J) \models \gamma$ . By Lemma 2, this is true if and only if  $\not\vdash_{\mathcal{B}} \gamma \Rightarrow \bigvee J$ .
- Pick a strengthening constraint  $\bigwedge J \in S$ , and let  $v \in S_-$  with  $v = c(\bigwedge J)$  be the corresponding negative sample. Moreover, assume that  $\gamma$  is consistent with  $\mathfrak{S}$  and, thus,  $\not\vdash_{\mathcal{B}} \bigwedge J \Rightarrow \gamma$ . By Lemma 2, this is true if and only if  $c(\bigwedge J) \not\models \gamma$ . Hence,  $v \not\models \gamma$ .  
Conversely, assume that  $\gamma$  is consistent with  $\mathcal{S}$ . Thus,  $v \not\models \gamma$ , which means  $c(\bigwedge J) \not\models \gamma$ . By Lemma 2, this is true if and only if  $\not\vdash_{\mathcal{B}} \bigwedge J \Rightarrow \gamma$ .
- Following the definition of implication, we split the proof into two cases, depending on whether  $\not\vdash_{\mathcal{B}} \bigwedge J \Rightarrow \gamma$  or  $\not\vdash_{\mathcal{B}} \gamma \Rightarrow \bigvee J$  (and  $v_1 \not\models \gamma$  or  $v_2 \models \gamma$  for the reserve direction). However, the proof in the former case is the same as the proof for strengthening constraints, while the proof of latter case is the same as the proof for weakening. Hence, combining both proofs immediately yields the claim.  $\square$

**ICE learners for Boolean formulas** The reduction above allows us to use any ICE learning algorithm in the literature that synthesizes positive Boolean formulas. As we mentioned earlier, we can add the negations of predicates as first-class predicates, and hence synthesize invariants over the class of all Boolean combinations of a finite set of predicates as well.

The problem of passive ICE learning for one round, synthesizing a formula that satisfies the ICE sample, can usually be achieved efficiently and in a variety of ways. However, the crucial aspect is not the complexity of learning in one round, but the *number* of rounds it takes to converge to an adequate invariant that proves the program correct. When the set  $\mathcal{P}$  of candidate predicates is large (hundreds in our experiments), since the number of Boolean formulas over  $\mathcal{P}$  is doubly exponential in  $n = |\mathcal{P}|$ , building an effective learner is not easy. However, there is one class of formulas that are particularly amenable to efficient ICE learning—learning *conjunctions of predicates over  $\mathcal{P}$* . In this case, there are ICE learning algorithms that promise learning the invariant (provided one exists expressible as a conjunct over  $\mathcal{P}$ ) in  $n + 1$  rounds. Note that this learning is essentially finding an invariant in a hypothesis class  $\mathcal{H}$  of size  $2^n$  in  $n + 1$  rounds.

HOUDINI [21] is such a learning algorithm for conjunctive formulas. Though it is typically seen as a particular way to synthesize invariants, it is a prime example of an ICE learner for conjuncts, as described in the work by Garg et al. [24]. In fact, Houdini is similar to the classical PAC learning algorithm for conjunctions [35], but extended to the ICE model by handling implication counterexamples. More precisely, given an ICE sample  $\mathcal{S} = (S_+, S_-, S_{\Rightarrow})$ , HOUDINI computes the largest conjunctive formula  $\varphi$  in terms of the number of Boolean variables occurring in  $\varphi$  (i.e., the semantically strongest conjunctive formula) that is consistent with  $\mathcal{S}$  in the following way. First, it computes the largest conjunction  $\varphi$  that is consistent with the positive examples (i.e.,  $v \models \varphi$  for all  $v \in S_+$ ); note that this conjunction is unique. Next, HOUDINI checks whether the implications are satisfied. If this is not the case, then we know for each non-satisfied implication

```

int A[], B[];
int N;    axiom (N > 0);
bool inImage(int i) { return true; }

procedure inverse ()
requires (∀x, y. 0 ≤ x < y < N ⇒ A[x] ≠ A[y]); // A is injective
requires (∀x. 0 ≤ x < N ∧ inImage(x) ⇒ (∃y. 0 ≤ y < N ∧ A[y] = x)); // A is surjective
ensures (∀x, y. 0 ≤ x < y < N ⇒ B[x] ≠ B[y]); // B is injective
{
  int i = 0;
  while (i < N)
  SynthesizeInv (∀x. 0 ≤ x < i ⇒ B[A[x]] = x); // b1
  {
    B[A[i]] = i;
    i = i + 1;
  }
  SynthesizeInv (∀x. 0 ≤ x < N ⇒ A[B[x]] = x, // b2
                ∀x. 0 ≤ x < N ∧ inImage(x) ⇒ A[B[x]] = x); // b3
  return;
}

```

**Fig. 2.** Synthesizing invariants for the program that constructs an inverse B of an injective, surjective function A [36].

$(v_1, v_2) \in S_{\Rightarrow}$  that  $v_2$  has to be classified positively because  $v_1$  belongs to every set that includes  $S_+$ . Hence, HOUDINI adds all such  $v_2$  to  $S_+$ , resulting in a new set  $S'_+$ . Subsequently, it constructs the largest conjunction  $\varphi'$  that is consistent with the positive examples in  $S'_+$  (i.e.,  $v \models \varphi'$  for all  $v \in S'_+$ ). HOUDINI repeats this procedure until it arrives at the largest conjunctive formula  $\varphi^*$  that is consistent with  $S_+$  and  $S_{\Rightarrow}$  (again, note that this set is unique). Finally, HOUDINI checks whether each negative example violates  $\varphi^*$  (i.e.,  $v \not\models \varphi^*$  for all  $v \in S_-$ ). If this is the case,  $\varphi^*$  is the largest conjunctive formula over  $\mathcal{B}$  that is consistent with  $\mathcal{S}$ ; otherwise, no consistent conjunctive formula exists. The time HOUDINI spends in each round is *polynomial* and, furthermore, when used in an iterative setting, is guaranteed to converge in at most  $n + 1$  rounds or report that no conjunctive invariant over  $\mathcal{P}$  exists. We use this ICE learner to build a CD-NPI learner for conjunctions.

### 2.3 An Illustrative Example

Figure 2 illustrates an example program of the verified software competition [36] that given an injective, surjective function A returns the inverse B of the function A. The post-condition of this program expresses that the function B is injective. To prove this program correct, one needs to specify adequate invariants at the loop header and before the return statement in the function *inverse* in the program. We wish to synthesize these invariants. For simplicity, let us assume we are provided a small set of predicates as building blocks of the invariants to synthesize— $b_1$  for the loop invariant and  $b_2, b_3$  for the invariant at the return statement. Our task, therefore, is to synthesize adequate invariants for this program over these predicates.<sup>6</sup>

Clearly, the verification conditions of this program are undecidable. In fact, the constant Boolean function *inImage* is crucially required to validate certain verification conditions in BOOGIE because it triggers appropriate quantifier instantiations in the surjectivity condition.

Now, suppose the learner conjectures the loop invariant  $\gamma_L = b_1$  and the invariant at the return statement  $\gamma_R = b_2 \wedge b_3$ . Moreover, suppose that the verification condition along the path from the loop exit to the return statement, though valid in the undecidable theory  $\mathcal{U}$  (cf. Figure 1), is not provable in the decidable theory  $\mathcal{D}$

<sup>6</sup> In general, one starts with a much larger set of candidate predicates that are automatically generated using program/specification-dependent heuristics.

(one that has instantiated quantifiers with ground terms). The  $\mathcal{D}$ -solver returns a model for the negation of the verification condition that captures this non-provability information. The verification engine gleans this model—it looks for the values assigned to the predicate variables in the model, and from this information constructs, in general, a CD-NPI constraint for the learner to learn from. For this particular verification, the verification engine extracts a pair of formulas  $(\eta, \chi)$  where  $\eta = b_1$  and  $\chi = b_2$ , and communicates this as an inductivity constraint to the learner. Intuitively, this constraint means that the verification condition obtained by substituting  $\gamma_L$  with  $\eta$  and  $\gamma_R$  with  $\chi$  is itself not provable. Hence, in subsequent rounds, the learner needs to conjecture only such invariants where  $\gamma_L$  is not weaker than  $\eta$  (i.e.,  $\not\vdash_{\mathcal{B}} b_1 \Rightarrow \gamma_L$ ) or  $\gamma_R$  is not stronger than  $\chi$  (i.e.,  $\not\vdash_{\mathcal{B}} \gamma_R \Rightarrow b_2$ ).

The learner works by reducing the CD-NPI passive learning problem to ICE learning over a sample over the given set of predicates. Concretely, the inductivity constraint  $(b_1, b_2)$  is reduced to an implication constraint  $((1, 0, 0), (1, 0, 1))$  in the ICE setting, where each datapoint in the ICE sample has values for the predicates  $b_1$ ,  $b_2$ , and  $b_3$ , respectively. In the next round, let us assume the learner conjectures the invariants  $\gamma_L = b_1$  and  $\gamma_R = b_3$ . Note these conjectures satisfy both the ICE constraints and the CD-NPI constraints. In this case, it turns out that the verification conditions along all program paths using these invariants can be proved valid in the theory  $\mathcal{D}$ . As a result, our invariant synthesis procedure terminates with  $\gamma_L$  and  $\gamma_R$  as adequate inductive invariants.

## 2.4 Main Result

To state the main result of this paper, let us assume that the set  $\mathcal{P}$  of predicates is finite. We comment on the case of infinitely many predicates below.

**Theorem 2.** *Assume a normal verification engine for a program  $P$  to be given. Moreover, let  $\mathcal{P}$  be a finite set of predicates over the variables in  $P$  and  $\mathcal{H}$  a hypothesis class consisting of positive Boolean combinations of predicates in  $\mathcal{P}$ . If there exists an annotation in  $\mathcal{H}$  that the verification engine can use to prove  $P$  correct, then the CD-NPI framework described in Section 2.1 is guaranteed to converge to such an annotation in finite time.*

*Proof (Proof of Theorem 2).* The proof proceeds in two steps. First, we show that a normal verification engine is *honest*, meaning that the non-provability information returned by such an engine does not rule out any adequate and provable annotation. Second, we show that any consistent learner (i.e., a learner that only produces consistent hypotheses), when paired with an honest verification engine, makes *progress* from one round to another. Finally, we combine both results to show that the framework eventually converges to an adequate and provable annotation.

*Honesty of the verification engine* We show honesty of the verification engine individually for each type of constraint by contradiction.

- Suppose that the verification replies to a candidate invariant  $\gamma$  proposed by the learner with a weakening constraint  $\chi$  because it could not prove the validity of the Hoare triple  $\{\alpha\}s\{\gamma\}$ . This effectively forces any future conjecture  $\gamma'$  to satisfy  $\not\vdash_{\mathcal{B}} \gamma' \Rightarrow \chi$ .  
Now, suppose that there exists an invariant  $\delta$  such that  $\vdash_{\mathcal{B}} \delta \Rightarrow \chi$  and the verification engine can prove the validity of  $\{\alpha\}s\{\delta\}$  (in other words, the adequate invariant  $\delta$  is ruled out by the weakening constraint  $\chi$ ). Due to the fact that the verification engine is normal (in particular, by contraposition of Part 1 of Definition 1), this implies that the verification engine can also prove the validity of  $\{\alpha\}s\{\chi\}$ . However, this is a contradiction to  $\chi$  being a weakening constraint.
- Suppose that the verification engine replies to a candidate invariant  $\gamma$  proposed by the learner with a strengthening constraint  $\eta$  because it could not prove the validity of the Hoare triple  $\{\gamma\}s\{\beta\}$ . This effectively forces any future conjecture  $\gamma$  to satisfy  $\not\vdash_{\mathcal{B}} \eta \Rightarrow \gamma$ .  
Now, suppose that there exists an invariant  $\delta$  such that  $\vdash_{\mathcal{B}} \eta \Rightarrow \delta$  and the verification engine can prove the validity of  $\{\delta\}s\{\beta\}$  (in other words, the adequate invariant  $\delta$  is ruled out by the weakening constraint  $\eta$ ). Due to the fact that the verification engine is normal (in particular, by contraposition of Part 2 of

Definition 1), this implies that the verification engine can also prove the validity of  $\{\eta\}s\{\beta\}$ . However, this is a contradiction to  $\eta$  being a strengthening constraint.

- Combining the arguments for weakening and strengthening constraints immediately results in a contradiction for the case of inductivity constraints as well.

*Progress of the learner* Now suppose that the learning algorithm is consistent, meaning that it always produces an annotation that is consistent with the current sample. Moreover, assume that the sample in iteration  $i \in \mathbb{N}$  is  $\mathfrak{S}_i$  and the learner produces the annotation  $\gamma_i$ . If  $\gamma_i$  is inadequate to prove the program correct, the verification engine returns a constraint  $c$ . The learner adds this constraint to the sample, obtaining the sample  $\mathfrak{S}_{i+1}$  of the next iteration.

Since verification with  $\gamma_i$  failed, which is witnessed by  $c$ , we know that  $\gamma_i$  is not consistent with  $c$ . The next conjecture  $\gamma_{i+1}$ , however, is guaranteed to be consistent with  $\mathfrak{S}_{i+1}$  (which contains  $c$ ) because the learner is consistent. Hence,  $\gamma_i$  and  $\gamma_{i+1}$  are semantically different. Using this argument repeatedly shows that each annotation  $\gamma_i$  that a consistent learner has produced is semantically different from any previous annotation  $\gamma_j$  for  $j < i$ .

*Convergence* We first make two observations.

1. The number of semantically different hypotheses in the hypothesis space  $\mathcal{H}$  is finite because the set  $\mathcal{P}$  is finite. Recall that  $\mathcal{H}$  is the class of all positive Boolean combinations of predicates in  $\mathcal{P}$ .
2. Due to the honesty of the verification engine, every annotation that the verification engine can use to prove the program correct is guaranteed to be consistent with any sample produced during the learning process.

Now, suppose that there exists an annotation that the verification engine can use to prove the program correct. Since the learner is consistent, all conjectures produced during the learning process are semantically different. Thus, the learner will at some point have exhausted all incorrect annotations in  $\mathcal{H}$  (due to Observation 1). By assumption, however, there exists at least one annotation that the verification engine can use to prove the program correct. Moreover, any such annotation is guaranteed to be consistent with the current sample (due to Observation 2). Thus, the annotation conjectured next is necessarily one that the verification engine can use to prove the program correct.  $\square$

Under certain, realistic assumptions on the CD-NPI learning algorithm, Theorem 2 remains true even if the number of predicates is infinite. An example of such an assumption is that the learning algorithm always conjectures a smallest consistent annotation with respect to some fixed total order on  $\mathcal{H}$ . In this case, one can show that such a learner will at some point have proposed all inadequate annotation up to the smallest annotation the verification engine can use to prove the program correct. It will then conjecture this annotation in the next iteration. We refer the reader to Löding, Madhusudan, and Neider [41] for details on further strategies that ensure convergence.

### 3 Application: Learning Invariants that Aid Natural Proofs for Heap Reasoning

We now develop an implementation of our learning framework for verification engines based on natural proofs for heap reasoning [50,48]. We first provide some background on the separation logic DRYAD and natural proofs, which is a sound but incomplete verification procedure. Then, we describe how to implement our verification framework using a natural proofs verification engine. In particular, we describe how to automatically generate suitable predicates for these programs, which serve as the building blocks of the invariants we seek to synthesize. Finally, we present an empirical evaluation of our implementation on an extensive set of standard algorithms on dynamic data structures, such as searching, inserting, or deleting items in lists and trees.

**Background: Natural Proofs and Dryad** DRYAD [50,48] is a dialect of separation logic that comes with a heaplet semantics and allows expressing second order properties such as pointing to a list (or list segment) using recursive functions and predicates. The syntax of DRYAD is a standard separation logic syntax with a few restrictions, such as disallowing negations inside recursive definitions and in sub-formulas connected by spatial conjunction (see [48] for more details about the DRYAD syntax). DRYAD is expressive enough to state a variety of data-structures (singly and doubly linked lists, sorted lists, binary search trees, AVL trees, maxheaps, treaps, etc.), recursive definitions over them that map to numbers (length, height, etc.), as well as data stored within the heap (the multiset of keys stored in lists, trees, etc.).

Natural proofs [50,48] is a sound but incomplete strategy for deciding satisfiability of DRYAD formulas. The first step of the natural proof verifier is to convert all predicates and functions in a DRYAD-annotated program to *classical logic*. This translation introduces *heaplets* (modeled as sets of locations) explicitly in the logic. Furthermore, it introduces assertions that demand that the accesses of each method are contained in the heaplet implicitly defined by its precondition (taking into account newly allocated or freed nodes), and that at the end of the program, the modified heaplet precisely matches the implicit heaplet defined by the post-condition.

The second step of the natural proof verifier is to perform *transformations* on the program and translate it to BOOGIE [21], an intermediate verification language that handles proof obligations using automatic theorem provers (typically SMT solvers). This transformation essentially performs three tasks: (a) it abstracts all recursive definitions on the heap using uninterpreted functions and introduces finite-depth unfoldings of recursive definitions at every place in the code where locations are dereferenced, (b) it models heaplets and other sets using a decidable theory of maps, and (c) it inserts *frame reasoning* explicitly in the code that allows the verifier to derive that certain properties continue to hold across a heap update (or function call) using the heaplet that is modified. The resulting BOOGIE program is a program with no recursive definitions, where all verification conditions are in decidable logics, and where the logic engine can return models when formulas are satisfiable. The program can be verified if supplied with correct inductive loop-invariants and adequate pre/post conditions.

The described procedure has been implemented in a fully automatic tool, called VCDRYAD. VCDRYAD extends VCC [13] and converts C programs annotated in DRYAD to BOOGIE programs via the natural proof transformations described above. It is important to note, however, that VCC introduces some quantification to define the memory model and semantics of C, but this does not typically derail decidable reasoning. We refer the reader to [50,48] for more details.

**Learning Heap Invariants** We have implemented a prototype of our CD-NPI framework over VCDRYAD and the BOOGIE program verifier. This prototype takes a C program annotated in DRYAD as input and uses VCDRYAD to convert it to a BOOGIE program. Then, it applies our transformation to the ICE learning framework and automatically generates a set  $\mathcal{P}$  of *predicates* (as described shortly), which serve as the basic building blocks of our invariants. Finally, it pairs the BOOGIE verifier with an invariant synthesis engine, HOUDINI in our case, to learn an inductive invariant. Note that after the VCDRYAD-transformation, BOOGIE satisfies the requirements on verification engines of our framework.

The set  $\mathcal{P}$  of predicates is generated from generic templates, shown in Figure 3, which are instantiated using all combinations of program variables that occur in the program being verified. The templates define a fairly exhaustive set of predicates, including

- properties of the store (equality of pointer variables, equality and inequalities between integer variables, etc.),
- shape properties (singly and doubly linked lists and list segments, sorted lists, trees, BST, AVL, treaps, etc.),
- and recursive definitions that map data structures to numbers (keys/data stored in a structure, lengths of lists and list segments, height of trees) involving arithmetic relationships and set relationships.

In addition, there are also predicates describing heaplets of various structures (with suffix *\_heaplet*), involving set operations, disjointness, and equalities. The structures and predicates are extensible, of course, to any recursive definition expressed in DRYAD.

$x, y \in \text{PointerVars}$      $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \text{PointerVars}^*$      $pf \in \text{PointerFields}$      $key, df \in \text{DataFields}$   
 $i, j \in \text{IntegerVars} \cup \{0, \text{IntMax}, \text{IntMin}\}$

$listshape(\mathbf{x}) := \text{LinkedList}(x_1) \mid \text{DoublyLinkedList}(x_1) \mid \text{SortedLinkedList}(x_1)$ $\quad \mid \text{LinkedListSeg}(x_1, x_2) \mid \text{DoublyLinkedListSeg}(x_1, x_2)$ $\quad \mid \text{SortedLinkedListSeg}(x_1, x_2)$		
$treeshape(\mathbf{x}) := \text{BST}(x) \mid \text{AVLtree}(x) \mid \text{Treap}(x)$		
$shape(\mathbf{x}) := listshape(\mathbf{x}) \mid treeshape(\mathbf{x})$		
$size(\mathbf{x}) := listshape\_length(\mathbf{x}) \mid treeshape\_height(\mathbf{x})$		
Category 1	$x = \text{nil}$ $x \neq \text{nil}$ $shape(\mathbf{x})$ $x \in shape\_heaplet(\mathbf{y})$ $x \notin shape\_heaplet(\mathbf{y})$ $shape\_heaplet(\mathbf{x}) \cap shape\_heaplet(\mathbf{y}) = \emptyset$	$x = y$ $x \neq y$ $x.pf = \text{nil}$ $x.pf \neq \text{nil}$ $x.pf = y$ $x.pf \neq y$
Category 2	$i \in shape\_key\_set(\mathbf{x})$ $i \notin shape\_key\_set(\mathbf{x})$ $shape\_key\_set(\mathbf{x}) \leq_{\text{set}} \{i\}$ $shape\_key\_set(\mathbf{x}) \geq_{\text{set}} \{i\}$ $shape\_key\_set(\mathbf{x}) \leq_{\text{set}} \{y.df\}$ $shape\_key\_set(\mathbf{x}) \geq_{\text{set}} \{y.df\}$ $shape\_key\_set(\mathbf{x}) = shape\_key\_set(\mathbf{y})$ $shape\_key\_set(\mathbf{x}) \leq_{\text{set}} shape\_key\_set(\mathbf{y})$ $shape\_key\_set(\mathbf{x}) \geq_{\text{set}} shape\_key\_set(\mathbf{y})$ $shape\_key\_set(\mathbf{x}) = shape\_key\_set(\mathbf{y})$ $\quad \cup shape\_key\_set(\mathbf{z})$	$x.df = i$ $x.df \neq i$ $x.df \leq i$ $x.df \geq i$ $x.df = y.df$ $x.df \neq y.df$ $x.df \leq y.df$ $x.df \geq y.df$
Category 3	$size(\mathbf{x}) = i - j$ $size(\mathbf{x}) - size(\mathbf{y}) = i$ $size(\mathbf{x}) - size(\mathbf{y}) = i - j$	$size(\mathbf{x}) = i$ $size(\mathbf{x}) \leq i$ $size(\mathbf{x}) \geq i$

**Fig. 3.** Templates for generating predicates. The operator  $\leq_{\text{set}}$  denotes comparison between integer sets, where  $A \leq_{\text{set}} B$  if and only if  $\forall x \in A. \forall y \in B. x \leq y$ . The operator  $\geq_{\text{set}}$  is similarly defined. Shape properties such as `LinkedList`, `AVLtree`, etc., are recursively defined in DRYAD, separately, and is extensible to any class of DRYAD defined shapes. Similarly, the definitions related to keys stored in a datastructure and the sizes of datastructures also stem from recursive definitions of them in DRYAD.

The predicates are grouped into three categories, roughly in increasing complexity. Category 1 predicates involve shape-related properties, Category 2 involves properties related to the keys stored in the data-structure, and Category 3 predicates involve size-predicates on data structures (lengths of lists and heights of trees). Given a program to verify and its annotations, we choose the category of predicates depending on whether the specification refers to shape only, shapes and keys, or shapes, keys, and sizes (choosing a category includes the predicates of lower category as well). Then, predicates are automatically generated by instantiating the

templates with all (combinations of) program variables. This approach allows for a fine-grained control over the predicates that are generated for a specific program and prevents the set of predicates from growing too large.

**Evaluation** We have evaluated our prototype on ten benchmark suits (82 routines in total) that contain standard algorithms on dynamic data structures, such as searching, inserting, or deleting items in lists and trees. These benchmarks were taken from the following sources: (1) GNU C Library(glibc) singly/sorted linked lists, (2) GNU C Library(glibc) doubly linked lists, (3) OpenBSD SysQueue, (4) GRASSHOPPER [49] singly linked lists, (5) GRASSHOPPER [49] doubly linked lists, (6) GRASSHOPPER [49] sorted linked lists, (7) VCDRYAD [48] sorted linked lists, (8) VCDRYAD [48] binary search trees, AVL trees, and treaps, (9) AFWP [32] singly/sorted linked lists, and (10) ExpressOS [42] MemoryRegion. The specifications for these programs are generally checks for their full functional correctness, such as preserving or altering shapes of data structures, inserting or deleting keys, filtering or finding elements, and sortedness of elements. The specifications hence involve separation logic with arithmetic as well as recursive definitions that compute numbers (like lengths and heights), data-aggregating recursive functions (such as multisets of keys stored in data-structures), and complex combinations of these properties (e.g., to specify binary search trees, AVL trees and treaps). All programs are annotated in DRYAD, and checking validity of the resulting verification conditions is undecidable.

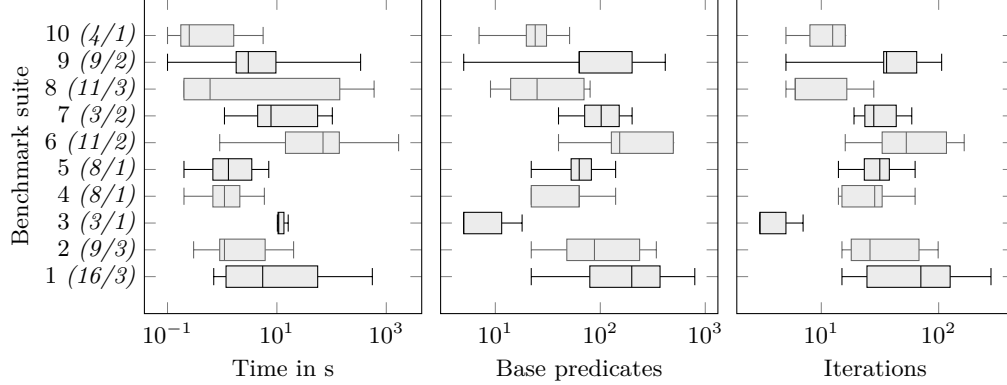
To create our benchmarks, we first picked all programs that contained iterative loops, *erased* the user-provided loop invariants, and used our framework to synthesize adequate inductive invariants (our tool can synthesize multiple invariants for a program). We also selected some programs that were purely recursive, where the contract for the function had been strengthened to make the verification succeed. We *weakened* these contracts to only state the specification (typically by removing formulas in the post-conditions of recursively-called functions) and introduced annotation holes instead. The goal was to synthesize strengthenings of these contracts that allow proving the program correct. We also chose five straight-line programs, deleted their post-conditions, and evaluated whether we can learn post-conditions for them. Since our conjunctive learner learns the strongest invariant expressible as a conjunct, we can use our framework to synthesize post-conditions as well.

After removing annotations from the benchmarks, we automatically inserted appropriate predicates over which to build invariants and contracts as described above. For all benchmark suits, conjunctions of these predicates were sufficient to prove the program correct.

*Experimental Results* We performed all experiments in a virtual machine running Ubuntu 16.04.1 on a single core of an Intel Core i7-7820 HK 2.9 GHz CPU with 2 GB memory. The box plots in Figure 4 summarize the results of this empirical evaluation aggregated by benchmark suite, specifically the time required to verify the programs, the number of base predicates, and the number iterations in the learning process (see Appendix A for full details). Each box in the diagrams shows the lower and upper quartile (left and right border of the box, respectively), the median (line within the box), as well as the minimum and maximum (left and right whisker, respectively).

Our prototype was successful in learning invariants and contracts for all 82 benchmarks. Moreover, the fact that the median time for a great majority of benchmarks suits is less than 10 s shows that our technique is extremely effective in finding inductive DRYAD invariants. We also observe that despite many examples having hundreds of base predicates, which in turn suggests a worst-case complexity of hundreds of iterations, the learner was able to learn with much fewer iterations and the number of predicates in the final invariant is small. This shows that the non-provability information provided by the natural proof engine provides much more information than what the worst-case suggests.

To the best of our knowledge, our prototype is the only tool currently able of fully automatically verifying this challenging benchmark set. We must emphasize, however, that there are subsets of our benchmarks that can be solved by reformulating verification in decidable fragments of separation logic studied—we refer the reader to the related work in Section 1 for a survey of such work. Our goal in this evaluation, however, is not to compete with other, mature tools on a subset of benchmarks, but to measure the efficacy of our proposed CD-NPI based invariant synthesis framework on the whole benchmark set.



**Fig. 4.** Experimental evaluation of our prototype. The numbers in italic brackets shows the total number of programs in the suite (first number) and the maximum predicate category used (second number).

## 4 Application: Learning Invariants in the Presence of Bounded Quantifier Instantiation

Software verification of numerous applications must deal with quantification. For instance, quantifiers are often needed for axiomatizing theories that are not already equipped with decision procedures, for specifying properties of unbounded data structures and dynamically allocated memory, as well as for defining recursive properties of programs. For instance, the power of two function can be defined recursively using quantifiers as

$$pow2(0) = 1 \text{ and } \forall n \in \mathbb{N}: n > 0 \Rightarrow pow2(n) = 2 \cdot pow2(n - 1).$$

Despite the fact that various important first-order theories are undecidable (e.g., the first-order theory of arithmetic with uninterpreted functions), modern SMT solvers implement a host of heuristics to cope with quantifier reasoning. Quantifier instantiation, including pattern-based quantifier instantiation (e.g., E-matching [17]) and model-based quantifier instantiation [26], are particularly effective heuristics in this context. The key idea of instantiation-based heuristics is to instantiate universally quantified formulas with a finite number of ground terms and then check for validity of the resulting quantifier-free formulas (whose theory needs to be decidable). The exact instantiation of ground terms varies from method to method, but most instantiation-based heuristics are necessarily incomplete in general due to the undecidability of the underlying decision problems.

We can apply invariant synthesis framework for verification engines that employ quantifier instantiation in the following way. Assume that  $\mathcal{U}$  is an undecidable first-order theory allowing uninterpreted functions and that  $\mathcal{D}$  is its decidable quantifier-free fragment. Then, quantifier instantiation can be seen as a transformation of a  $\mathcal{U}$ -formula  $\varphi$  (potentially containing quantifiers) into a  $\mathcal{D}$ -formula  $approx(\varphi)$  in which all existential quantifiers have been eliminated (e.g., using skolemization) and all universal quantifiers have been replaced by finite conjunctions over ground terms.<sup>7</sup> This means that if the  $\mathcal{D}$ -formula  $approx(\varphi)$  is valid, then the  $\mathcal{U}$ -formula  $\varphi$  is valid as well. On the other hand, if  $approx(\varphi)$  is not valid, one cannot deduce the validity of  $\varphi$ . However, a  $\mathcal{D}$ -model of  $approx(\varphi)$  can be used to derive non-provability information as described in Section 2.1.

We have implemented our learning framework for synthesizing invariants based on bounded quantifier instantiation. Our prototype is based on BOOGIE/Z3 as the verification engine and uses HOUDINI to learn conjunctive invariants. In the remainder of this section, we present empirical results of this implementation on benchmarks taken from competitions and verified systems such as IronFleet [29].

<sup>7</sup> Quantifier instantiation is usually performed iteratively, but we here abstract away from this fact.



**Evaluation** We collected a benchmarks suite of twelve programs, which we obtained by simplifying programs found in IronFleet [29] (provably correct distributed systems), VSCOMP (Verified Software Competition) benchmarks [36], ExpressOS [42] (a secure operating system for mobile devices), and sparse matrix multiplication programs [8]. In these programs, quantifiers are used in specifying recursively defined predicates such as  $\text{power}(n, m)$  and  $\text{sum}(n)$ , and various array properties such as minimum/maximum elements, existence of specific elements, no duplicate elements, permutations of array elements, relations between two arrays, periodic properties of array elements, and bijective (injective and surjective) maps. The specifications hence are undecidable and fall outside of the decidable array property fragment [7]. In particular, the array specifications involve strict comparison ( $<$ ) between universally quantified index variables, array accesses in the index guard, nested array accesses (e.g.,  $a_1[a_2[i]]$ ), arithmetic expressions over universally quantified index variables, and alternation of universal and existential quantifiers.

From this benchmark suite, we erased the user-defined loop invariants and used our framework to find adequate inductive invariants. We generated a set of predicates that serve as the building blocks of our invariants. To this end, we used the pre-/post-conditions of the program being verified as templates from which the actual predicates are generated; the templates are instantiated using all combinations of program variables that occur in the program. We also generated predicates for octagonal constraints, (i.e., relations between two integer variables of the form,  $\pm x \pm y \leq c$ ). For a few programs, we also generated the octagonal predicates over array access expressions that appear in the program.

*Experimental Results* We performed all experiments in a virtual machine running Ubuntu 16.04.1 on a single core of an Intel Core i7-7820 HK 2.9 GHz CPU with 2 GB memory. The results of these experiments are listed in Table 1.

**Table 1.** Experimental results of the quantifier instantiation benchmarks. The column  $|\mathcal{P}|$  refer to the number of candidate predicates, the column  $\# \text{ Iterations}$  to the number of iterations of the teacher and learner, and the column  $|\text{Inv}|$  to the number of predicates in the inferred invariant.

Program	$ \mathcal{P} $	$\# \text{ Iterations}$	$ \text{Inv} $	Time in s
inverse	414	126	73	9.04
power2	109	55	34	2.10
powerN	160	60	31	13.52
recordArraySplit	1264	49	51	57.46
recordArrayUnzip	222	17	25	0.84
removeDuplicates	280	67	86	4.43
setFind	492	74	136	2.76
setInsert	556	73	188	6.70
sparseMatrixGen	816	278	90	22.07
sparseMatrixMul	768	313	91	14.49
sum	128	40	22	1.02
sumMax	192	61	45	4.31

As can be seen from the table, our framework is effective in finding inductive invariants that result in proving the programs correct (with an average of less than a minute per routine). Despite having hundreds of candidate predicates in many examples, which in turn suggests a worst-case complexity of hundreds of rounds, the learner was able to learn with much fewer rounds. Again, the non-provability information provided by the verification engine provides much more information than the worst-case suggests.

## 5 Conclusion

We have presented learning-based framework for invariant synthesis in the presence of sound but incomplete verification engines. To prove that our technique is effective in practice, we have implemented our framework

for two types of specifications: an expressive and undecidable dialect of separation logic called DRYAD for specifying heap properties and specifications involving universal quantification. In both cases, our prototype turned out to be extremely effective in learning inductive invariants and pre/post-conditions. In particular, the benchmark suite for DRYAD-annotated programs is extremely challenging, containing an extensive list of standard algorithms on dynamic data structures, and we are not aware of any other technique that can handle this benchmark suite.

Several future research directions are interesting. First, the framework we have developed is based on CEGIS where the invariant synthesizer synthesizes invariants using non-provability information but does not directly work on the program’s structure. It would be interesting to extend white-box invariant generation techniques such as interpolation/IC3/PDR, working using  $\mathcal{D}$  (or  $\mathcal{B}$ ) abstractions of the program directly in order to synthesize invariants for them. Second, in the NPI learning framework we have put forth, it would be interesting to change the underlying logic of communication  $\mathcal{B}$  to a richer logic, say the theory of arithmetic and uninterpreted functions. The challenge here would be to extract non-provability information from the models to the richer theory, and pairing them with synthesis engines that synthesize expressions against constraints in  $B$ . Finally, we think invariant learning should also include *experience* gained in verifying other programs in the past, both manually and automatically. A learning algorithm that combines logic-based synthesis with experience gained from repositories of verified programs can be more effective.

## References

1. Albarghouthi, A., Berdine, J., Cook, B., Kincaid, Z.: Spatial interpolants. In: ESOP 2015 (part of ETAPS 2015). NCS, vol. 9032, pp. 634–660. Springer (2015)
2. Ball, T., Majumdar, R., Millstein, T.D., Rajamani, S.K.: Automatic predicate abstraction of C programs. In: PLDI, 2001. pp. 203–213 (2001)
3. Barnett, M., Chang, B.E., DeLine, R., Jacobs, B., Leino, K.R.M.: Boogie: A modular reusable verifier for object-oriented programs. In: FMCO 2005. LNCS, vol. 4111, pp. 364–387. Springer (2005)
4. Betts, A., Chong, N., Donaldson, A.F., Qadeer, S., Thomson, P.: Gpuverify: a verifier for GPU kernels. In: OOPSLA 2012. pp. 113–132. ACM (2012)
5. Beyer, D., Keremoglu, M.E.: Cpathchecker: A tool for configurable software verification. In: CAV, 2011. pp. 184–190 (2011)
6. Bradley, A.R.: Sat-based model checking without unrolling. In: VMCAI 2011. LNCS, vol. 6538, pp. 70–87. Springer (2011)
7. Bradley, A.R., Manna, Z., Sipma, H.B.: What’s decidable about arrays? In: VMCAI, 2006. pp. 427–442. VMCAI’06, Springer-Verlag (2006)
8. Buluç, A., Fineman, J.T., Frigo, M., Gilbert, J.R., Leiserson, C.E.: Parallel sparse matrix-vector and matrix-transpose-vector multiplication using compressed sparse blocks. In: SPAA, 2009. pp. 233–244. SPAA ’09, ACM (2009)
9. Calcagno, C., Distefano, D., O’Hearn, P.W., Yang, H.: Compositional shape analysis by means of bi-abduction. J. ACM 58(6), 26 (2011)
10. Chen, E.Y., Chen, S., Qadeer, S., Wang, R.: Securing multiparty online services via certification of symbolic transactions. In: SP 2015. pp. 833–849. IEEE Computer Society (2015)
11. Chin, W.N., David, C., Nguyen, H.H., Qin, S.: Automated verification of shape, size and bag properties via user-defined predicates in separation logic. Sci. Comput. Program. 77(9), 1006–1036 (2012)
12. Chu, D., Jaffar, J., Trinh, M.: Automatic induction proofs of data-structures in imperative programs. In: PLDI, 2015. pp. 457–466. ACM (2015)
13. Cohen, E., Dahlweid, M., Hillebrand, M.A., Leinenbach, D., Moskal, M., Santen, T., Schulte, W., Tobies, S.: VCC: A practical system for verifying concurrent C. In: TPHOLs 2009. LNCS, vol. 5674, pp. 23–42. Springer (2009)
14. Cohen, E., Paul, W.J., Schmaltz, S.: Theory of multi core hypervisor verification. In: SOFSEM 2013. Lecture Notes in Computer Science, vol. 7741, pp. 1–27. Springer (2013)
15. Colón, M., Sankaranarayanan, S., Sipma, H.: Linear invariant generation using non-linear constraint solving. In: CAV 2003. LNCS, vol. 2725, pp. 420–432. Springer (2003)
16. Cousot, P., Cousot, R.: Abstract interpretation: A unified lattice model for static analysis of programs by construction or approximation of fixpoints. In: POPL 1977. pp. 238–252. ACM Press (1977)
17. Detlefs, D., Nelson, G., Saxe, J.B.: Simplify: a theorem prover for program checking. J. ACM 52(3), 365–473 (2005)

18. Dillig, I., Dillig, T., Li, B., McMillan, K.L.: Inductive invariant generation via abductive inference. In: OOPSLA 2013. pp. 443–456 (2013)
19. Een, N., Mishchenko, A., Brayton, R.: Efficient implementation of property directed reachability. In: FMCAD, 2011. pp. 125–134. FMCAD '11, FMCAD Inc (2011)
20. Ernst, M.D., Czeisler, A., Griswold, W.G., Notkin, D.: Quickly detecting relevant program invariants. In: ICSE 2000. pp. 449–458. ACM Press (2000)
21. Flanagan, C., Leino, K.R.M.: Houdini, an annotation assistant for `esc/java`. In: FME 2001. LNCS, vol. 2021, pp. 500–517. Springer (2001)
22. Floyd, R.W.: Assigning Meanings to Programs. In: Schwartz, J.T. (ed.) Proceedings of a Symposium on Applied Mathematics. Mathematical Aspects of Computer Science, vol. 19, pp. 19–31. American Mathematical Society (1967)
23. Garg, P., Löding, C., Madhusudan, P., Neider, D.: Learning universally quantified invariants of linear data structures. In: CAV 2013. LNCS, vol. 8044, pp. 813–829. Springer (2013)
24. Garg, P., Löding, C., Madhusudan, P., Neider, D.: ICE: A robust framework for learning invariants. In: CAV 2014. LNCS, vol. 8559, pp. 69–87. Springer (2014)
25. Garg, P., Madhusudan, P., Neider, D., Roth, D.: Learning Invariants using Decision Trees and Implication Counterexamples. In: POPL, 2016. pp. 499–512 (2016)
26. Ge, Y., de Moura, L.M.: Complete instantiation for quantified formulas in satisfiability modulo theories. In: CAV 2009. pp. 306–320 (2009)
27. Gulwani, S., Srivastava, S., Venkatesan, R.: Program analysis as constraint solving. In: PLDI 2008. pp. 281–292. ACM (2008)
28. Gupta, A., Rybalchenko, A.: Invgen: An efficient invariant generator. In: CAV 2009. LNCS, vol. 5643, pp. 634–640. Springer (2009)
29. Hawblitzel, C., Howell, J., Kapritsos, M., Lorch, J.R., Parno, B., Roberts, M.L., Setty, S.T.V., Zill, B.: Ironfleet: proving practical distributed systems correct. In: SOSP 2015, 2015. pp. 1–17 (2015)
30. Hawblitzel, C., Howell, J., Lorch, J.R., Narayan, A., Parno, B., Zhang, D., Zill, B.: Ironclad apps: End-to-end security via automated full-system verification. In: OSDI 2014. pp. 165–181. USENIX Association (2014)
31. Hoare, C.A.R.: An axiomatic basis for computer programming. *Commun. ACM* 12(10), 576–580 (1969)
32. Itzhaky, S., Banerjee, A., Immerman, N., Nanevski, A., Sagiv, M.: Effectively-propositional reasoning about reachability in linked data structures. In: CAV 2013. Lecture Notes in Computer Science, vol. 8044, pp. 756–772. Springer (2013)
33. Itzhaky, S., Bjørner, N., Reps, T.W., Sagiv, M., Thakur, A.V.: Property-directed shape analysis. In: CAV, 2014. pp. 35–51 (2014)
34. Karbyshev, A., Bjørner, N., Itzhaky, S., Rinetzky, N., Shoham, S.: Property-directed inference of universal invariants or proving their absence. In: CAV 2015. pp. 583–602 (2015)
35. Kearns, M.J., Vazirani, U.V.: An Introduction to Computational Learning Theory. MIT Press (1994)
36. Klebanov, V., Müller, P., Shankar, N., Leavens, G.T., Wüstholtz, V., Alkassar, E., Arthan, R., Bronish, D., Chapman, R., Cohen, E., Hillebrand, M., Jacobs, B., Leino, K.R.M., Monahan, R., Piessens, F., Polikarpova, N., Ridge, T., Smans, J., Tobies, S., Tuerk, T., Ulbrich, M., Weiß, B.: The 1st Verified Software Competition: Experience report. In: FM, 2011. LNCS, vol. 6664. Springer (2011)
37. Krishna, S., Puhersch, C., Wies, T.: Learning invariants using decision trees. CoRR abs/1501.04725 (2015)
38. Lal, A., Qadeer, S., Lahiri, S.K.: A solver for reachability modulo theories. In: CAV 2012. LNCS, vol. 7358, pp. 427–443. Springer (2012)
39. Le, Q.L., Gherghina, C., Qin, S., Chin, W.: Shape analysis via second-order bi-abduction. In: CAV 2014. pp. 52–68 (2014)
40. Leino, K.R.M.: Dafny: An automatic program verifier for functional correctness. In: LPAR 2010. LNCS, vol. 6355, pp. 348–370. Springer (2010)
41. Löding, C., Madhusudan, P., Neider, D.: Abstract learning frameworks for synthesis. In: TACAS 2016. LNCS, vol. 9636, pp. 167–185. Springer (2016)
42. Mai, H., Pek, E., Xue, H., King, S.T., Madhusudan, P.: Verifying security invariants in expressos. In: ASPLOS 2013. pp. 293–304. ACM (2013)
43. McMillan, K.L.: Interpolation and SAT-Based model checking. In: CAV 2003. LNCS, vol. 2725, pp. 1–13. Springer (2003)
44. de Moura, L.M., Bjørner, N.: Efficient e-matching for SMT solvers. In: CADE, 2007. pp. 183–198 (2007)
45. de Moura, L.M., Bjørner, N.: Z3: an efficient SMT solver. In: TACAS 2008. LNCS, vol. 4963, pp. 337–340. Springer (2008)

46. Padhi, S., Sharma, R., Millstein, T.D.: Data-driven precondition inference with learned features. In: PLDI, 2016. pp. 42–56 (2016)
47. Pavlinovic, Z., Lal, A., Sharma, R.: Inferring annotations for device drivers from verification histories. In: ASE, 2016. pp. 450–460 (2016)
48. Pek, E., Qiu, X., Madhusudan, P.: Natural proofs for data structure manipulation in C using separation logic. In: PLDI 2014. p. 46. ACM (2014)
49. Piskac, R., Wies, T., Zufferey, D.: Automating separation logic using SMT. In: CAV 2013. Lecture Notes in Computer Science, vol. 8044, pp. 773–789. Springer (2013)
50. Qiu, X., Garg, P., Stefanescu, A., Madhusudan, P.: Natural proofs for structure, data, and separation. In: PLDI 2013. pp. 231–242. ACM (2013)
51. Sagiv, M., Reps, T., Wilhelm, R.: Parametric shape analysis via 3-valued logic. In: POPL, 1999. pp. 105–118. POPL '99, ACM (1999)
52. Sharma, R., Aiken, A.: From invariant checking to invariant inference using randomized search. In: CAV 2014. LNCS, vol. 8559, pp. 88–105. Springer (2014)
53. Sharma, R., Gupta, S., Hariharan, B., Aiken, A., Liang, P., Nori, A.V.: A data driven approach for algebraic loop invariants. In: ESOP 2013. LNCS, vol. 7792, pp. 574–592. Springer (2013)
54. Sharma, R., Gupta, S., Hariharan, B., Aiken, A., Nori, A.V.: Verification as learning geometric concepts. In: SAS 2013. LNCS, vol. 7935, pp. 388–411. Springer (2013)
55. Sharma, R., Nori, A.V., Aiken, A.: Interpolants as classifiers. In: CAV 2012. LNCS, vol. 7358, pp. 71–87. Springer (2012)
56. Yang, J., Hawblitzel, C.: Safe to the last instruction: automated verification of a type-safe operating system. In: PLDI 2010. pp. 99–110 (2010)
57. Zhu, H., Nori, A.V., Jagannathan, S.: Learning refinement types. In: ICFP, 2015. pp. 400–411. ICFP 2015, ACM (2015)

## A Detailed Results of the Heap Invariants Benchmarks

Table 2: Experimental results of the heap invariants benchmarks. The column  $|\mathcal{P}|$  refer to the number of candidate predicates, the column *Cat.* corresponds to the category of predicates used, the column *# Iterations* to the number of iterations of the teacher and learner, and the column  $|\text{Inv}|$  to the number of predicates in the inferred invariant. A  $\dagger$  indicates contract strengthening, while a  $*$  indicates post condition learning.

(1) GNU C Library(glibc) Singly and Sorted Linked-List

Program	$ \mathcal{P} $	Cat.	# Iterations	$ \text{Inv} $	Time in s
g_slist_copy	368	2	123	101	55
g_slist_find	48	2	18	9	0.8
g_slist_free	22	1	15	1	1.2
g_slist_index	237	3	68	57	6.3
g_slist_insert	464	2	160	50	219.1
g_slist_insert_before	795	2	279	114	556.1
g_slist_insert_sorted	520	2	193	135	210.6
g_slist_last	32	2	19	8	0.7
g_slist_length	54	3	20	12	0.9
g_slist_nth	88	3	26	17	1.1
g_slist_nth_data	342	3	99	62	9.2
g_slist_position	162	3	32	18	2.7
g_slist_remove	140	1	73	28	4.7
g_slist_remove_all	380	2	132	15	57.7
g_slist_remove_link	325	1	85	57	14.3
g_slist_reverse	117	2	58	6	4.5

(2) GNU C Library(glibc) Doubly Linked-List

Program	$ \mathcal{P} $	Cat.	# Iterations	$ \text{Inv} $	Time in s
g_list_find	48	2	18	9	0.8
g_list_free	22	1	15	1	0.9
g_list_index	237	3	68	57	6.1
g_list_last	22	1	15	6	0.3
g_list_length	88	3	24	16	1.1
g_list_nth	88	3	26	17	1.1
g_list_nth_data	342	3	99	62	9.4
g_list_position	162	3	32	18	2.9
g_list_reverse	320	2	84	2	20.2

(3) OpenBSD SysQueue

Program	$ \mathcal{P} $	Cat.	# Iterations	$ \text{Inv} $	Time in s
squeue_insert_head*	5	1	3	2	10.8
squeue_insert_tail*	5	1	3	3	16.1
squeue_remove_head*	18	1	7	5	10.2

Table 2: continued

(4) GRASShopper [49] Singly Linked-List

Program	$ \mathcal{P} $	Cat.	# Iterations	Inv	Time in s
sl_concat	63	1	26	15	0.9
sl_copy	63	1	32	12	2.7
sl_dispose	22	1	14	1	0.7
sl_filter	140	1	63	9	5.9
sl_insert	63	1	31	19	1.3
sl_remove	22	1	15	6	0.6
sl_reverse	63	1	36	4	1.9
sl_traverse	22	1	15	4	0.2

(5) GRASShopper [49] Doubly Linked-List

Program	$ \mathcal{P} $	Cat.	# Iterations	Inv	Time in s
dl_concat	63	1	26	15	0.7
dl_copy	63	1	32	12	3.1
dl_dispose	140	1	44	4	7.1
dl_filter	140	1	63	9	4.6
dl_insert	63	1	31	19	0.9
dl_remove	22	1	15	6	0.4
dl_reverse	63	1	36	4	1.7
dl_traverse	22	1	14	4	0.2

(6) GRASShopper [49] Sorted Linked-List

Program	$ \mathcal{P} $	Cat.	# Iterations	Inv	Time in s
sls_concat	153	2	38	27	11
sls_copy	496	2	144	94	1679.7
sls_dispose	40	2	16	3	1.6
sls_double_all	496	2	106	102	118.8
sls_filter	496	2	127	30	158.3
sls_insert	153	2	53	29	17.7
sls_merge	416	2	63	29	327.0
sls_remove	496	2	165	119	69.7
sls_reverse	102	2	28	13	21.6
sls_split	153	2	53	29	98.1
sls_traverse	40	2	18	9	0.9

(7) VCDryad [48] Sorted Linked-List

Program	$ \mathcal{P} $	Cat.	# Iterations	Inv	Time in s
find_last_sorted	40	2	19	11	1.1
reverse_sorted	102	2	28	13	7.8
sorted_insert_iter	201	2	59	48	103.2

Table 2: continued

## (8) VCDryad [48] Trees

Program	$ \mathcal{P} $	Cat.	# Iterations	Inv	Time in s
avl-delete-rec <sup>†</sup>	72	3	16	5	449.1
avl-find-smallest <sup>†</sup>	19	3	5	11	0.2
avl-insert-rec <sup>†</sup>	72	3	23	14	102.2
bst-delete-rec <sup>†</sup>	68	2	16	11	180
bst-find-rec <sup>†</sup>	23	2	6	9	0.5
bst-insert-rec <sup>†</sup>	68	2	28	16	64.8
traverse-inorder <sup>†</sup>	9	3	6	3	0.2
traverse-posttorder <sup>†</sup>	9	3	6	3	0.2
traverse-preorder <sup>†</sup>	9	3	6	3	0.2
treap-delete-rec <sup>†</sup>	80	3	17	13	599.0
treap-find-rec <sup>†</sup>	25	3	6	11	0.6

## (9) AFWP [32] Singly and Sorted Linked-List

Program	$ \mathcal{P} $	Cat.	# Iterations	Inv	Time in s
SLL-create	5	1	5	1	0.1
SLL-delete-all	22	1	14	1	5.3
SLL-delete	265	1	106	47	9.6
SLL-filter	63	1	34	9	2.1
SLL-find	140	1	53	45	3
SLL-insert	201	2	65	26	45.6
SLL-last	63	1	34	9	1.2
SLL-merge	416	2	71	46	339.3
SLL-reverse	63	1	36	4	1.8

## (10) ExpressOS [42] MemoryRegion

Program	$ \mathcal{P} $	Cat.	# Iterations	Inv	Time in s
memory_region_find	24	1	16	2	0.2
memory_region_init*	7	1	5	4	0.1
memory_region_insert	51	1	16	3	0.3
split_memory_region*	24	1	9	6	5.6