

# Semantic Foundations for the Static Analysis of Program Revisions

DAKOTA BRYAN, University of Colorado Boulder, USA

## 1 Introduction

During the software development process, programmers create pull requests (PRs) to modify the behavior of a system. These PRs are subsequently reviewed to determine if the changes introduced to the code produce the intended consequences. Often during these reviews, informal notions of the semantics of the change are used. For example, “in the old version the value of  $x$  could be 10, but the new version does not allow that.” Or “the only states that we’re removed are one’s where the value of  $x$  is 17”. Or even, “in the states when the values of  $x$  is the same in both the new and old version, the values of  $y$  are different”. In this work, we present a concrete semantics that can describe the semantic effect of program revision, and an abstraction interpretation over it which enables reasoning akin the examples above.

There are many program analyses that enable reasoning about two programs; however, we claim that they fall short by not capturing the mental model developers hold of diffs and not bringing that model to the forefront of the analysis. One existing technique, “Abstract Semantic Differencing for Numerical Programs” [9] is designed specifically for program differences, and has a concept of “added”, “removed”, and “related” states. However, we claim their version of determining if a variable value is removed does not capture all ways that a programmer may consider a value to be removed. It is fundamentally based on a state-to-state comparison (akin to relational Hoare logic) which does not allow for the set of “added” states to have a different cardinality as the set of “removed” states.

We center our analysis around a characterization of “old” and “new” behaviors as added, removed, or related. In order to enable the possibility for the added and removed set of behaviors to have a different cardinality, we must form this partition via a state-to-set of state comparison. The checks needed to form this partition of possible states is not always obvious enough to encode with RHL. Existing analyses which can form a similar partition do not provide a natural way to write assertions over them, or shape the analysis around them.

### 1.1 Contributions

Our main contributions can be summarized as follows:

- We first give a concrete semantics that describes the semantic effect of program revision.
- Then, we define an abstract semantics over this, which allows for proofs about program revisions.

## 2 Overview

In this section, we work through an example program revision and prove differential properties about the revision. We see how our natural intuitions about program differences lead to the formalisms and generalizations given in section 3. Particularly, we see that a natural way to describe program differences leads to the concrete semantics of our “diff machine” given in section 3.2. Furthermore, we see that this “diff machine” must be parametrized by some notion of state equivalence, which generalizes to the *state correspondence* given in section 3.1. Finally, in an attempt

---

Author’s Contact Information: Dakota Bryan, University of Colorado Boulder, Boulder, Colorado, USA, dakota.bryan@colorado.edu.

Listing 1. (Sign Program)

```

50
51   x := nonDet()
52   if (x > 0)
53       sgn := 1
54   else
55       sgn := -1
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98

```

Listing 2. (Patched Sign Program)

```

59   x := nonDet()
60   if (x > 0)
61       sgn := 1
62   else
63       sgn := -1
64   if (x = 0)
65       sgn := 0
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98

```

Fig. 1. An example revision to a program that computes the sign of a number

to prove safety properties of the reachable “diff machine” states, we notice several intricacies which are crucial to the design of our parameterized abstract interpreter (section 3.5).

We will consider imperative programs written in the IMP language with non-deterministic assignment (the grammar can be found in appendix A).

Consider the program and corresponding revision given in figure 1. We will let program variables subscripted with  $-$  denote variables in the “old” program and  $+$  subscripts for “new” program variables. For example,  $x_-$  refers to  $x$  in program 1. We will also let  $\Sigma_-$  denote the reachable states of program 1 and  $\Sigma_+$  denote the set of reachable states for program 2. We see that the unrevised program (1) computes the sign of  $x_-$  with zero considered negative. The revised program (2) adds the functionality that the sign of zero is zero.

## 2.1 Describing the Semantic Effect of Program Revision

Before proving properties about program revisions, we need a way to represent revisions; then, we can prove properties over that representation. Therefore, our first task is to describe the semantic effect of a revision (this will generalize to our concrete semantics). As we saw in the section 1, software developers can view a program revision in terms of the behavior which it removed, added, and kept the same. For IMP programs, one measure of behavior is the set of reachable states. Thus, a first attempt of describing the semantic effect of a revision is characterizing reachable states of the two programs as removed, related, or added. To form the removed set, **rm**, we could check if some reachable state of the old program is not reachable in the new program, and if so, that state would be removed. This set will be considered the behavior that the revision removed.

Now, considering our example and this approach, we would have  $\sigma_- := x \mapsto 0, \text{sgn} \mapsto -1$  as a removed state. However, one may view the effect of the revision as adding the functionality where the sign of a number can now be zero. With this view, the revision removes no behavior. In general, the set of removed states should need not have the same cardinality as the set of added states. For example, a revision that removes bugs but adds no functionality is naturally modeled with **add** (the set of “added” behavior) empty and **rm** filled with the buggy states.

In order to address this, we can consider a  $\sigma_-$  to be removed only if its value of sgn is not a possible value of sgn in  $\Sigma_+$  ( $\mathbf{rm} := \{\sigma_- \in \Sigma_- \mid \forall \sigma_+ \in \Sigma_+. \sigma_-(\text{sgn}) \neq \sigma_+(\text{sgn})\}$ ). Thus, whereas before we built  $\mathbf{rm}$  by checking  $\sigma_- \neq \sigma_+$  for every  $\sigma_+ \in \Sigma_+$ , we now check some projection of the two states. So, with this updated definition, we have  $\mathbf{rm} := \{\}$ . We can compute  $\mathbf{add}$  the same way to achieve  $\mathbf{add} := \{x \mapsto 0, \text{sgn} \mapsto 0\}$ .

Now, what about behavior that the revision kept the same? We first notice that if there is some  $\sigma_- \notin \mathbf{rm}$ , then there must be some  $\sigma_+ \in \Sigma_+$  where  $\sigma_-(\text{sgn}) = \sigma_+(\text{sgn})$ . This implies that each reachable old state is either removed or *related* to some reachable state for the new program, when the relation used to compare them is the state projection check. This “projection check” will henceforth be referred to as the *state correspondence*. Now, for checking properties over the semantic effect of this revision, it is natural to check properties about the behavior that is not removed and not added. So, as a first guess of defining the “related”, say  $\mathbf{rel}$ , part of our “diff machine” (*i.e.* semantic description of the effect of the revision, *i.e.* concrete semantics), we let  $\mathbf{rel} := (\Sigma_- \setminus \mathbf{rm}) \times (\Sigma_+ \setminus \mathbf{add})$ .

This construction certainly captures the behaviors that were not removed or added, but we would have the following state pair in  $\mathbf{rm}$ :  $(x_- \mapsto 5, \text{sgn}_- \mapsto 1), (x_+ \mapsto -5, \text{sgn}_+ \mapsto -1)$ . This state pair isn’t exactly *related*. Instead, we want every  $(\sigma_-, \sigma_+) \in \mathbf{rel}$  to be related, or satisfy the state correspondence:  $\text{sgn}_- = \text{sgn}_+$ . This is now a little more natural; we could prove that in all related states, if  $x_- > 0$ , then  $x_+ > 0$ . This alternate formulation of  $\mathbf{rel}$  is the same as  $(\Sigma_- \setminus \mathbf{rm}) \times (\Sigma_+ \setminus \mathbf{add})$  with the state correspondence added as an implication antecedent to every assertion over it. This new formulation still characterizes the entire reachable state space such that  $\forall \sigma_- \in \Sigma_-. \sigma_- \in \mathbf{rm} \vee \exists \sigma_+. (\sigma_-, \sigma_+) \in \mathbf{rel}$ . We “leave nothing out”.

## 2.2 Reasoning About Program Revisions

Equipped with the description of the semantic effect of a program revision ( $\mathbf{rm}, \mathbf{rel}, \mathbf{add}$ ), we turn our attention to proving safety properties over this description. We will continue with the program revision shown in figure 1, along with the particular  $\mathbf{rm}, \mathbf{rel}, \mathbf{add}$  for that program revision explained in section 2.1.

Consider  $\varphi_{\mathbf{rm}}$  as a safety property over  $\mathbf{rm}$ ; so,  $\varphi_{\mathbf{rm}}$  relies on old program variables ( $\varphi_{\mathbf{rm}}(\sigma_-)$ ). And  $\varphi_{\mathbf{rel}}$  as a safety property over  $\mathbf{rel}$ , which is a relation property over two states ( $\varphi_{\mathbf{rel}}(\sigma_-, \sigma_+)$ ). Now, we have  $(\varphi_{\mathbf{rm}}, \varphi_{\mathbf{rel}}, \varphi_{\mathbf{add}})$  as a safety property over  $\mathbf{rm}, \mathbf{rel}, \mathbf{add}$ . We notice that  $\mathbf{rm}, \mathbf{rel}, \mathbf{add}$  is a certain filtered partition of  $\Sigma_+ \times \Sigma_+$ ; so, we will attempt to use relational Hoare logic to prove  $(\varphi_{\mathbf{rm}}, \varphi_{\mathbf{rel}}, \varphi_{\mathbf{add}})$ .

Recall the semantics of relational Hoare logic, shown below.

$$\models \{\varphi\} s_- \pm s_+ \{\varphi'\} \text{ if and only if } \\ [(\sigma_-, \sigma_+) \models \varphi \text{ and } \sigma_- \vdash s_- \Downarrow \sigma'_- \text{ and } \sigma_+ \vdash s_+ \Downarrow \sigma'_+] \text{ implies } (\sigma'_-, \sigma'_+) \models \varphi'$$

Let  $\varphi_{\mathbf{add}} := x_+ = 0$ ,  $\varphi_{\mathbf{rel}} := x_- > 0 \implies x_+ > 0$ ,  $s_- :=$  program 1, and  $s_+ :=$  program 2. Let’s first attempt to prove  $\varphi_{\mathbf{rel}}$ ; recall,  $(\sigma_-, \sigma_+) \in \mathbf{rel}$  if and only if  $\sigma_-(\text{sgn}) = \sigma_+(\text{sgn})$ . So, we can build the post condition of our relational Hoare triple by first filtering state pairs which satisfy our state correspondence, then checking if those satisfy  $\varphi_{\mathbf{rel}}$ .  $\varphi' := \text{sgn}_- = \text{sgn}_+ \implies \varphi_{\mathbf{rel}}$ . Great,  $\{\top\} s_- \pm s_+ \{\varphi'\}$  is derivable so we can prove our assertion over  $\mathbf{rel}$ , indicating that safety properties about behaviors which a revision preserves can be proven using relational Hoare logic.

Now, let’s attempt to prove  $\varphi_{\mathbf{add}}$  over  $\mathbf{add}$ . Recall,  $\sigma_+ \in \mathbf{add}$  if and only if, for all  $\sigma_- \in \Sigma_-$ ,  $\sigma_+(\text{sgn}) \neq \sigma_-(\text{sgn})$ . How can we encode that the state correspondence does not hold for any  $\sigma_-$  to use as our filtering antecedent? Post conditions in relational Hoare logic are interpreted over all reachable state pairs. So if we wrote  $\text{sgn} \neq \text{sgn}$  in the post condition, it would be satisfied by  $(\sigma_- = x_- \mapsto 5, -\text{sgn} \mapsto 1), (\sigma_+ = x_+ \mapsto -5, \text{sgn} \mapsto -1)$ . This is not what we want. In fact, there is

no way to encode our property into relational Hoare logic without oracle power to determine the set of  $\text{sgn}$  values in the old program. This is because our property is *hyper*, for one new state, we need to compare it against the *set* of reachable old states. Relational Hoare logic (or Hoare logic) assertions must be over *states*, not *sets of states*.

The fact that related properties can be proven in relational Hoare logic but not added or removed properties is an indication as to why much attention has been placed on proving *equivalence* of program revisions, but not *difference*. Furthermore, to ensure that  $|\text{rm}| \neq |\text{add}|$  it is insufficient to compare on the state level. If the two states were found to be not related, then the old one would go in **rm** and the new one in **add**, guaranteeing the sets to be equal cardinality.

Now, let's take a slightly different approach. We can use Hoare logic to find a formula that represents all possible sign values in the old program. Consider the derivable Hoare logic triple,  $\{\top\} s_- \{\varphi_-\}$  where  $\varphi_- = (x_- > 0 \implies \text{sgn}_- = 1) \wedge (x_- \leq 0 \implies \text{sgn}_- = -1)$ .  $\varphi_-$  contains all the information we need for the possible values of  $\text{sgn}_-$ ; so, we can use this in combination with a formula representing reachable states of the new program (say  $\varphi_+$ ) to prove  $\varphi_{\text{add}}$ . We just need some formula that checks if  $\text{sgn}_- \neq \text{sgn}_+$  holds for all  $\sigma_- \in \Sigma_-$ , given  $\varphi_-$  that represents  $\Sigma_-$ . We can universally quantify over variables modified by program 1 to achieve  $\forall x_-, \text{sgn}_-. \varphi_-$  which eliminates  $\varphi_-$ 's dependence of the variables in program 1. This has the effect of transforming  $\varphi_-$  into a *global* property (from a state property) that must hold for every possible assignment of those variables; which is exactly what we desire.

Now, we can add  $\varphi_-$  as the antecedent to  $\text{sgn}_- \neq \text{sgn}_+$  in a formula that universally quantifies over  $x_-$  and  $\text{sgn}_-$  to achieve a formula which represents all old states such that  $\text{sgn}_- \neq \text{sgn}_+$  holds for possible values of  $\text{sgn}$  in the old program, shown below.

$$\forall x_-, \text{sgn}_-. (\varphi_- \implies \text{sgn}_- \neq \text{sgn}_+)$$

Just as the quantification removes  $\varphi_-$ 's dependence on  $\sigma_-$ , it also the dependence on  $\sigma_-$  for  $\text{sgn}_- \neq \text{sgn}_+$ . Thus, the entire formula is interpreted over  $\sigma_+$ , as desired. This formula implies  $\varphi_r := 1 \neq \text{sgn}_+ \wedge -1 \neq \text{sgn}_+$ , which describes all old states that satisfy our removed state condition.

Now, given the derivable Hoare triple  $\{\top\} s_+ \{\varphi_+\}$  where  $\varphi_+ = (x_+ > 0 \implies \text{sgn}_+ = 1) \wedge (x_+ = 0 \implies \text{sgn}_+ = 0) \wedge (x_+ < 0 \implies \text{sgn}_+ = -1)$ , we can create  $\varphi_+ \wedge \varphi_r$ . This has the effect of filtering states described by  $\varphi_+$  (reachable old program states) to only those that satisfy  $\varphi_r$  (states that satisfy our removed state condition). And notice  $\varphi_+ \wedge \varphi_r \implies x_+ = 0$ , which was our desired safety assertion over **add**. Perfect, we have proven our **add** property; however, we notice the following issue. Since Hoare logic is an over approximating program logic, our old program triple,  $\{\top\} s_- \{\varphi_-\}$  where  $\varphi_- = (x_- > 0 \implies \text{sgn}_- = 1) \wedge (x_- \leq 0 \implies \text{sgn}_- = -1)$ , implies the triple  $\{\top\} s_- \{\varphi_- \vee \text{sgn}_- = 0\}$ . Then,  $\forall x_-, \text{sgn}_-. ((\varphi_- \vee \text{sgn}_- = 0) \implies \text{sgn}_- \neq \text{sgn}_+)$  implies  $\varphi_r := 1 \neq \text{sgn}_+ \wedge -1 \neq \text{sgn}_+ \wedge 0 \neq \text{sgn}_-$ . Now, we add the  $\varphi_+$  conjunct  $(\varphi_+ \wedge \varphi_r)$ , but this implies  $\perp$ .  $\perp$  is certainly not a sound over approximation of **add**.

Let's go back to our definition of added states, but consider approximations,  $\hat{\text{add}} := \{\sigma_+ \in \hat{\Sigma}_+ \mid \forall \sigma_- \in \hat{\Sigma}_-. \sigma_+ \langle \text{sgn} \rangle \neq \sigma_- \langle \text{sgn} \rangle\}$ . We want **add** to over approximate **add**. If  $\hat{\Sigma}_+$  over approximates  $\Sigma_+$  and  $\hat{\Sigma}_-$  is exact, then we can only add states to **add**. There might be some  $\sigma_+ \in \hat{\Sigma}_+$  but  $\notin \Sigma_+$  in which the state correspondence does not hold for all reachable old states, thus growing **add**. Now, if  $\hat{\Sigma}_+$  is exact, but  $\hat{\Sigma}_-$  is an over approximation, we can only remove states from **add**. Now, there may be some  $\sigma_- \in \hat{\Sigma}_-$  but  $\notin \Sigma_-$  in which some  $\sigma_+$  has an  $\text{sgn}_+$  that equals  $\sigma_- \langle \text{sgn} \rangle$ . Then, this  $\sigma_+$  would not be in **add**, but would be in **add**. Therefore, to soundly over approximate **add**, we require a representation of old states that under approximates all reachable states. Then, we may soundly over approximate **add**.

197 In order to under approximate reachable states, we use incorrectness logic. Recall the semantics,  
 198 show below.

$$\begin{aligned} 199 \quad & \models [\varphi] s [\varphi'] \text{ if and only if} \\ 200 \quad & \sigma' \models \varphi' \text{ implies } (\sigma \vdash s \Downarrow \sigma' \text{ and } \sigma \models \varphi) \end{aligned}$$

202 All free variables before the “implies” are implicitly universally quantified over, and all free variables  
 203 after the “implies” are implicitly existentially quantified over. Thus,  $\models [\varphi] s [\varphi']$  indicates that all  
 204 state satisfying  $\varphi'$  are reachable. This is precisely what we desire.  
 205

206 Now, we can derive the following incorrectness logic triple,  $[\top] s_- [\varphi_-]$  where  $\varphi_- = (x_- >$   
 207  $0 \implies \text{sgn}_- = 1) \wedge (x_- \leq 0 \implies \text{sgn}_- = -1)$ . We cannot soundly add disjunctions or drop  
 208 conjunctions in incorrectness logic. Then, our formula representing an over approximation of old  
 209 states that in which the state correspondence does not hold for all new states (under approximate)  
 210 is  $\forall x_-, \text{sgn}_-. \varphi_- \implies \text{sgn}_- \neq \text{sgn}_+$ . This implies  $\varphi_r := 1 \neq \text{sgn}_+ \wedge -1 \neq \text{sgn}_+$ . Then, finally,  
 211 creating the conjunction with our over approximate  $\varphi_+$  yields a formula which implies  $x_+ = 0$ .

212 In this section, we saw how partitioning reachable states of the “old” and “new” programs of a  
 213 revision into the sets **rm**, **rel**, **add** based on a state correspondence is a natural way to represent  
 214 the semantic effect of program revisions. Then, we noticed that in order to prove safety properties  
 215 over these sets, we require both an under and over approximating abstract semantics. Section 3  
 216 generalizes and formalizes these constructions. Finally, we used Hoare logic and Incorrectness logic  
 217 to prove a property over **add** and saw how relational Hoare logic was sufficient to prove **rel** safety  
 218 properties. Section 4 instantiates the parameterized abstract interpreter given in section 3 with the  
 219 two programs logics.

### 220 3 Parameterized Abstract Interpreter

221 In this section, we build up to a parameterized abstract interpreter for safety properties of removed,  
 222 related, and added behaviors caused by a program revision. We do not fix an underlying program-  
 223 ming language; programs of the underlying language are denoted as  $s$ . We write  $s_- \pm s_+$  to signify  
 224 a program  $s_-$  that was updated to  $s_+$  by a revision (sometimes referred to as a *bivisual* program).  
 225 If we consider a diff,  $\Delta$ , to be the changes made by a programmer to  $s_-$ , then  $s_+$  may be viewed as  
 226 the diff applied to the old program ( $s_+ = \Delta \circ s_-$ ). We write  $\sigma$  to represent a behavior of a program  $s$   
 227 that we will henceforth call *state*. A program state could be a store, trace, or some other notion of  
 228 program behavior, but we will leave it uninstantiated.  
 229

#### 230 3.1 State Correspondence

231 As we saw in the section 1, software developers can view a program revision in terms of the  
 232 behavior (or states) which it removed, added, and kept the same. So, the set of *removed* states  
 233 (denoted **rm**) is a subset of the set of all possible states of an old program ( $\Sigma_-$ ). Thus, any  $\sigma_- \in \Sigma_-$   
 234 is either  $\in \text{rm}$  or is *related* to some  $\sigma_+ \in \Sigma_+$ . The set of related states (**rel**) will be tuples of  $(\sigma_-, \sigma_+)$ ,  
 235 and if some  $\sigma_- \notin \text{rm}$ , then there will exist some  $(\sigma'_-, \sigma'_+) \in \text{rel}$  where  $\sigma_- = \sigma'_-$ . In fact, for all  
 236  $\sigma_i \in \Sigma_+$  that is related to  $\sigma_-$ ,  $(\sigma_-, \sigma_i)$  is in **rel**. The set of added states (**add**) includes the states  
 237 that were added by the revision, so states in **add** will not have some  $\sigma_- \in \Sigma_-$  that it is related to.  
 238 Together, the three sets (**rm**, **rel**, **add**) contain all possible program states of both programs and  
 239 characterize a revision.

240 Now, we see that to determine if a state is removed, we need an equivalence relation on states,  
 241 which we will call the *state correspondence* and use  $\sim$  to denote it. Given a state correspondence,  $\sim$ ,  
 242  $\sigma_-$  is removed if and only if for all  $\sigma_+ \in \Sigma_+$ ,  $\sigma_- \not\sim \sigma_+$  (where  $\sigma_- \not\sim \sigma_+$  holds if and only if  $\sigma_- \sim \sigma_+$   
 243 does not hold). This is quite natural, in order to check if some old behavior is removed, you must  
 244

check it against all new behaviors. Since we will be doing this for-all check quite a bit, we reserve the symbol,  $\neq$  for it, shown below.

**DEFINITION 1.** *The state-lifted state correspondence ( $\neq$ )*

$\sigma \neq \Sigma'$  if and only if for all  $\sigma' \in \Sigma'$ ,  $\sigma \not\sim \sigma'$

$\sigma \simeq \Sigma'$  holds if and only if  $\sigma \neq \Sigma'$  does not hold.

We write  $\simeq$  for the relation,  $\neq$ , not holding, which can be alternatively defined as:  $\sigma \simeq \Sigma'$  if and only if there exists some  $\sigma' \in \Sigma'$  where  $\sigma \sim \sigma'$ .

The state correspondence need not be the same relation when comparing old states to new states versus new states to old states. This is allowed in our framework.  $\neq^{\text{rm}}$  indicates that the state correspondence used to determine if an old state is removed is  $\sim^{\text{rm}}$ . Likewise for  $\neq^{\text{add}}$ ; if  $\sim$  and  $\neq$  is used, then it is assumed that  $\sim$  is the same for determining added and removed states.

We may now give definitions for **rm**, **rel**, **add**

**DEFINITION 2.** *The set of removed states (**rm**) based on a set of old states ( $\Sigma_-$ ), a set of new states ( $\Sigma_+$ ), and a state correspondence ( $\neq^{\text{rm}}$ ).*

$\text{rm} := \{\sigma_- \in \Sigma_- \mid \sigma_- \neq^{\text{rm}} \Sigma_+\}$

**DEFINITION 3.** *The set of related states (**rel**) based on a set of old states ( $\Sigma_-$ ), a set of new states ( $\Sigma_+$ ), and state correspondences ( $\neq^{\text{rm}}, \neq^{\text{add}}$ ).*

$\text{rel} := \{(\sigma_-, \sigma_+) \in \Sigma_- \times \Sigma_+ \mid \sigma_- \sim^{\text{rm}} \sigma_+ \text{ or } \sigma_- \sim^{\text{add}} \sigma_+\}$

The set of added states (**add**) is defined as expected. **rel** can be equivalently defined as  $\{(\sigma_-, \sigma_+) \in \Sigma_- \times \Sigma_+ \mid \sigma_- \sim^{\text{rm}} \sigma_+\} \cup \{(\sigma_-, \sigma_+) \in \Sigma_- \times \Sigma_+ \mid \sigma_- \sim^{\text{add}} \sigma_+\}$ . The left side represents all  $\sigma_- \notin \text{rm}$  paired with the new states they are related to. The right side represents all  $\sigma_+ \notin \text{add}$  which are related to some  $\sigma_-$  paired with its related old states. Now, we can see how this framework generalizes the concrete semantics of “all reachable pairs” frameworks like relational Hoare logic. If we let  $\sim$  hold universally, then **rel** is the set of all reachable pairs.

### 3.2 Concrete Semantics

Now we have a way to partition reachable states into the representation we are trying to prove safety properties over. However, we have not internalized the computation of reachable states; thus, we do not have a “machine” to abstract. If we had an abstraction or program logic that was expressive enough to soundly encode an over-approximation of **rm**, **rel**, **add**, we could use that and would not need to build the partition into the concrete semantics. However, we feel that the exercise of building the partition into the semantics is necessary to illuminate the intricacies of over-approximating these sets.

We will define our concrete semantics with the following judgment form:  $\Sigma_-, \Sigma_+ \vdash s_- \pm s_+ \Downarrow \text{rm}, \text{rel}, \text{add}$ . Its “input”, or context, is sets of initial new and old states. Then, it computes our three sets, **rm**, **rel**, **add**, which can be seen as the “value” type. It will also be dependent on state correspondences. This looks more like a collecting semantics than concrete, so it may be tempting to define something like  $\sigma_-, \sigma_+ \vdash s_- \pm s_+ \Downarrow (\sigma'_-, (\sigma''_-, \sigma'_+), \sigma'_+)$ . Then, define a collecting semantics over this to lift each component to states and abstract those states. However, recall from section 2.1 that state comparisons do not give us enough information to determine if a reachable old state is in **rm** or **rel**; we require sets. Therefore, after stepping a  $\sigma_-$  to  $\sigma'_-$ , we would not be able to tell where to put it in right side of our concrete semantics relation.

It also may be tempting to define both sides as  $\Sigma_-, \Sigma_+$ , since that would give us sufficient information to compute **rm**, **rel**, **add**. But remember, we are giving the semantics of our analysis that is an over-approximation of **rm**, **rel**, **add**, so we want our semantics to compute **rm**, **rel**, **add**.

Otherwise, we would need to define a projection from to  $\Sigma_-$ ,  $\Sigma_+$  to **rm**, **rel**, **add**, but this confuses and complicates our abstraction and soundness condition. Furthermore, **rm**, **rel**, **add** on the left side would be unnatural since with no programs, there is no added or removed states. This restricts the compositionality of the semantics, but a version with the **rm**, **rel**, **add** triple on both sides was investigated, and still did not allow for compositionality. This leads us to our concrete semantics.

**DEFINITION 4.** *The concrete semantics for our parameterized abstract interpreter. Parameterized by state correspondences ( $\sim^{\text{rm}}$ ,  $\sim^{\text{add}}$ ), an underlying programming language, a state representation of it, and its semantics ( $\sigma \vdash s \Downarrow^* \sigma'$ ). We take the natural collecting lift of these semantics ( $\Sigma \vdash s \Downarrow^* \Sigma'$ ).*

$$\frac{\begin{array}{c} \Sigma_- \vdash s_- \Downarrow^* \Sigma'_- \quad \Sigma_+ \vdash s_+ \Downarrow^* \Sigma'_+ \\ \mathbf{rm} = \{\sigma'_- \in \Sigma'_- \mid \sigma'_- \not\sim^{\text{rm}} \Sigma'_+\} \\ \mathbf{add} = \{\sigma'_+ \in \Sigma'_+ \mid \sigma'_+ \not\sim^{\text{add}} \Sigma'_-\} \\ \mathbf{rel} = \{(\sigma'_-, \sigma'_+) \in \Sigma'_- \times \Sigma'_+ \mid \sigma'_- \sim^{\text{rm}} \sigma'_+ \text{ or } \sigma'_- \sim^{\text{add}} \sigma'_+\} \end{array}}{\Sigma_-, \Sigma_+ \vdash s_- \pm s_+ \Downarrow \mathbf{rm}, \mathbf{rel}, \mathbf{add}}$$

We see that to evaluate from  $\Sigma_-$ ,  $\Sigma_+$  to **rm**, **rel**, **add** based on a biversional program, we step the initial states then compute the the sets via the state correspondence. Now, we can see that **rm**, **rel**, **add** is a particular view of the biversional state space which requires the whole space in order to be computed. Notice that we are assuming that base programming language is imperative, meaning that its step relation is one in which the context type is the same as the value type. Allowing a underlying program semantics where the context type does not equal the value type is possible and fairly straightforward, but left as future work.

### 3.3 Collecting Semantics

Given our concrete semantics, we may now define a collecting semantics and abstract domains. First, let's take a slight detour to discuss the abstract interpretation recipe we will use.

**DEFINITION 5.** *A general abstract interpretation framework*

- (1) *Concrete semantics:  $\sigma \vdash s \Downarrow v$*
- (2) *Collecting semantics:  $\Sigma \vdash s \Downarrow^{\{\}} V$*
- (3) *Abstract semantics:  $\hat{\sigma} \vdash s \Downarrow^{\#} \hat{v}$*

Where  $\Sigma : \mathcal{P}(\sigma)$ ,  $V : \mathcal{P}(v)$  and  $\gamma(\hat{\sigma}) = \{\sigma \mid \sigma \models \hat{\sigma}\}$ ,  $\gamma(\hat{v}) = \{v \mid v \models \hat{v}\}$ . Often, for imperative languages,  $\sigma = v$ .

In order to define an abstract semantics, we first define a concrete and collecting semantics. The essence of this recipe is that our collecting semantics defines the best possible abstract interpreter. That is, the soundness of the abstraction is defined w.r.t. the collecting semantics. Furthermore, if we had a sound and complete abstract interpreter, it would simply be the collecting semantics. Thus, the collecting semantics could lose precision from simply collecting all possible concrete states. Another important aspect of this recipe is that the collected state,  $\Sigma$ , is a collection of  $\sigma$ , and our abstract state,  $\hat{\sigma}$ , concretizes to a collected state. This allows for  $\gamma$  to be defined as all  $\sigma$  that models a  $\hat{\sigma}$ .

Given the semantics outlined in definition 5, we can define soundness of a typical over-approximating abstract interpreter as follows.

**DEFINITION 6.** *Soundness of an over-approximating abstract interpretation framework defined with definition 5.*

$\hat{\sigma} \vdash s \Downarrow^{\#} \hat{v}$  and  $\gamma(\hat{\sigma}) \vdash s \Downarrow^{\{\}} V$  implies  $V \subseteq \gamma(\hat{v})$ .

We will write  $\Downarrow^{\uparrow}$  for an over-approximating abstract semantics.

Our recipe does not enforce the abstract interpreter to be over-approximating. Below, we define the soundness of an abstraction that under-approximates the collecting semantics for bug-finding or incorrectness reasoning.

**DEFINITION 7.** Soundness of an under-approximating abstract interpretation framework defined with definition 5.

$\hat{s} \vdash s \Downarrow^\# \hat{v}$  and  $\gamma(\hat{s}) \vdash s \Downarrow^{\{\}} V$  implies  $\gamma(\hat{v}) \subseteq V$ .

We will write  $\Downarrow$  for an under-approximating abstract semantics.

Now, we can follow this recipe to define the collecting semantics. Our collected values,  $V$ , must be collection of  $v$  which is **rm**, **rel**, **add**. The same holds for our collected states. In order to lose no precision, we will define the typical set-lifted semantics.

**DEFINITION 8.** The collecting semantics of our parameterized abstract interpreter.  $\sigma \vdash s_- \pm s_+ \Downarrow v$  is the concrete semantics of our abstract interpreter given in definition 4.

$\Sigma ::= \circ \mid \Sigma, (\mathbf{rm}, \mathbf{rel}, \mathbf{add})$

$$\frac{V = \{v \mid \sigma \in \Sigma \text{ and } \sigma \vdash s_- \pm s_+ \Downarrow v\}}{\Sigma \vdash s_- \pm s_+ \Downarrow^{\{\}} V}$$

Our collecting semantics collect sets of **rm**, **rel**, **add** triples, so our abstract value,  $\hat{v}$ , must concretize to a set of **rm**, **rel**, **add** triples. However, each biversional program and state correspondence will produce a singular **rm**, **rel**, **add**, and this triple is what safety properties will be checked against. Thus, it is natural for  $\hat{v}$  to be an over-approximation of the three sets. Then, its concretization would be a **rm**, **rel**, **add** triple.

### 3.4 First Abstraction

Since this set-of-sets is not what we want, our first abstraction will be taking the set-of-sets back down to a set (*i.e.* lowering our collecting semantics to concrete semantics).

**DEFINITION 9.** Our first abstraction.  $\hat{\sigma} := (\Sigma_-, \Sigma_+)$ ,  $\hat{v} := (\mathbf{rm}, \mathbf{rel}, \mathbf{add})$ .

Abstract evaluation

$$\frac{(\Sigma_-, \Sigma_+) \vdash s_- \pm s_+ \Downarrow \mathbf{rm}, \mathbf{rel}, \mathbf{add}}{(\Sigma_-, \Sigma_+) \vdash s_- \pm s_+ \Downarrow^\# \mathbf{rm}, \mathbf{rel}, \mathbf{add}}$$

Value abstraction function.

$$\alpha : \mathcal{P}(V) \rightarrow V$$

$$\alpha(V, (\mathbf{rm}, \mathbf{rel}, \mathbf{add})) := (\mathbf{rm} \cup \alpha(V)_1, \mathbf{rel} \cup \alpha(V)_2, \mathbf{add} \cup \alpha(V)_3)$$

$$\alpha(\circ) := (\{\}, \{\}, \{\})$$

State abstraction function.

$$\alpha : \mathcal{P}(\Sigma) \rightarrow \Sigma$$

$$\alpha(\Sigma, (\Sigma_-, \Sigma_+)) := (\mathbf{rm} \cup \alpha(\Sigma)_1, \mathbf{rel} \cup \alpha(\Sigma)_2)$$

$$\alpha(\circ) := (\{\}, \{\}, \{\})$$

Element-wise state abstraction function ( $\beta$ ) is the identity function.

393      *Value concretization function.*

$$\begin{aligned} 394 \quad & \mathbf{rm}, \mathbf{rel}, \mathbf{add} \models \mathbf{rm}', \mathbf{rel}', \mathbf{add}' \text{ if and only if} \\ 395 \quad & \mathbf{rm} \subseteq \mathbf{rm}' \text{ and } \mathbf{rel} \subseteq \mathbf{rel}' \text{ and } \mathbf{add} \subseteq \mathbf{add}' \\ 396 \quad & \gamma(\hat{v}) := \{v \mid v \models \hat{v}\} \end{aligned}$$

398      *State concretization function.*

$$\begin{aligned} 400 \quad & \Sigma_-, \Sigma_+ \models \Sigma'_-, \Sigma'_+ \text{ if and only if} \\ 401 \quad & \Sigma_- \subseteq \Sigma_- \text{ and } \Sigma_+ \subseteq \Sigma_+ \\ 402 \quad & \gamma(\hat{\sigma}) := \{\sigma \mid \sigma \models \hat{\sigma}\} \end{aligned}$$

403      Now, we will show that in many cases, we can use a much more precise value concretization  
404      function that makes this abstraction sound and complete. First, notice that “running” the abstract  
405      interpreter, will start by abstracting a single concrete state,  $\beta(\sigma)$  (in this case,  $\sigma = \Sigma_-, \Sigma_+$ ). In all of  
406      these cases, we can refine our abstraction to be sound and complete.  
407

408      Our concrete semantics form a total and well-defined function:  $\lambda(\sigma, s_- \pm s_+) . v$  such that  $\sigma \vdash s_- \pm s_+ \Downarrow v$  where  $\Downarrow$  is the concrete evaluation from definition 4. Thus, for any  $\sigma = \Sigma_-, \Sigma_+$ ,  
409       $\{\sigma\} \vdash s_- \pm s_+ \Downarrow^{\{\}} V$  where  $\Downarrow^{\{\}}$  is given in definition 8,  $V$  will have cardinality one. That is,  
410       $V = \{(\mathbf{rm}, \mathbf{rel}, \mathbf{add})\}$  where **rm**, **rel**, **add** is the value produced by concrete semantics for  $\sigma$ . This  
411      implies that  $\beta(\sigma) \vdash s_- \pm s_+ \Downarrow^{\#} \hat{v}$  will produce a  $\hat{v}$  which is exactly all possible concrete values. Then,  
412      we can use alter the value concretization function for the abstraction given in 9 to achieve a sound  
413      and complete abstraction.  
414

415      DEFINITION 10. *An alteration to definition 9 that achieves a sound and complete abstraction. The*  
416      *value concretization function is redefined as follows in the case when  $\hat{\sigma} = \beta(\sigma)$ .*

$$\begin{aligned} 417 \quad & \mathbf{rm}, \mathbf{rel}, \mathbf{add} \models \mathbf{rm}', \mathbf{rel}', \mathbf{add}' \text{ if and only if} \\ 418 \quad & \mathbf{rm} = \mathbf{rm}' \text{ and } \mathbf{rel} = \mathbf{rel}' \text{ and } \mathbf{add} = \mathbf{add}' \\ 419 \quad & \gamma(\hat{v}) := \{v \mid v \models \hat{v}\} \end{aligned}$$

420      Anytime  $\beta(\sigma)$  is used as the context for  $\Downarrow^{\#}$ , this redefined  $\Downarrow^{\#}$  is used.  
421

422      Regardless of the existence of the alteration to  $\gamma$  presented in definition 10, after this initial  
423      abstraction, we are left with the task of abstracting  $\hat{\sigma}, \hat{v}$  which is the same as abstracting  $\sigma, v$ . Furthermore,  
424      with the special case of  $\gamma$  when  $\hat{\sigma} = \beta(\sigma)$ , we can ensure that our initial abstraction loses no  
425      precision. Thus, completeness can be determined w.r.t. any additional abstractions of  $\sigma = (\Sigma_-, \Sigma_+)$   
426      and  $v = (\mathbf{rm}, \mathbf{rel}, \mathbf{add})$ .  
427

### 428      3.5 Parameterized Abstract Semantics

429      Here, we give a parameterized abstract interpreter which can be instantiated with different abstract  
430      domains for the underlying program state space,  $\Sigma$ . As we saw in section 3.4, we are left with the  
431      task of abstracting our concrete semantics defined in definition 4. Viewing our entire landscape of  
432      semantics, this is actually an abstraction of the abstraction given in definition 9.  
433

434      From section 2.2, recall that in order to soundly over approximate **add**, we need an over approximating  
435      abstraction of reachable new states with an under approximating abstraction of reachable  
436      old states. We will write,  $\Downarrow^{\uparrow}$  for an over approximating abstract semantics in a domain  $\mathcal{D}$  (and  
437       $\Downarrow^{\downarrow}$  for under). We let elements in  $\mathcal{D}$  be denoted as  $\hat{\sigma}$ . Abstract state elements,  $\hat{\sigma}$ , concretize to a  
438      set of states. Therefore, to abstract our concrete semantics context,  $\Sigma_-, \Sigma_+$ , we will have a pair of  
439      abstract state elements:  $\hat{\sigma}_-, \hat{\sigma}_+$ . These elements are in the same domain that will be used in the  
440      parameterized abstract semantics for the underlying programs. Furthermore, the abstraction of **rm**  
441

and **add** will also be in this domain since they must concretize to a set of states. Our abstraction of **rel** should concretize to a set of pairs of states, so it will be in the domain  $\mathcal{D} \times \mathcal{D}$ .

Now, we see our parameterized abstract interpreter is parameterized by a single domain,  $\mathcal{D}$ . This is akin to the fact that our concrete semantics (definition 4) was also parameterized by a state representation of the underlying programs. The only aspect of our concrete semantics left to abstract is the state correspondence check, which has some intricacies discussed below. We found that the choice of how to abstract this check is dependent on the domain  $\mathcal{D}$ ; so in our abstract interpreter, we leave the this check as a required operation on the domain.

**DEFINITION 11.** *Parameterized abstract semantics. Parameters are a programming language whose programs are denoted  $s$ , a domain  $\mathcal{D}$  whose elements are abstractions of concrete states of  $s$ , and an over and under approximating abstract semantics on  $\mathcal{D}$  as  $s$ . Furthermore,  $\mathcal{D}$  must be a lattice that implements  $\hat{\neq}$ . In doing that,  $\mathcal{D}$  will also implement  $\hat{\sim}$ .*

$$\frac{\hat{\sigma}_- \vdash s_- \Downarrow^\uparrow \hat{\sigma}'_- \quad \hat{\sigma}_+ \vdash s_+ \Downarrow^\uparrow \hat{\sigma}'_+ \quad \hat{\sigma}_- \vdash s_- \Downarrow^\downarrow \hat{\sigma}''_- \quad \hat{\sigma}_+ \vdash s_+ \Downarrow^\downarrow \hat{\sigma}''_+}{\hat{\sigma}'_- \hat{\neq}^{rm} \hat{\sigma}''_+ \sqsubseteq \hat{\mathbf{rm}} \quad \hat{\sigma}'_+ \hat{\neq}^{add} \hat{\sigma}''_- \sqsubseteq \hat{\mathbf{add}} \quad [(\hat{\sigma}'_- \hat{\sim}^{rm} \hat{\sigma}'_+) \sqcup (\hat{\sigma}'_- \hat{\sim}^{add} \hat{\sigma}'_+)] \sqsubseteq \hat{\mathbf{rel}}}$$

$$\hat{\sigma}_-, \hat{\sigma}_+ \vdash s_- \pm s_+ \Downarrow^\# \hat{\mathbf{rm}}, \hat{\mathbf{add}}, \hat{\mathbf{rel}}$$

First, we notice the use of  $\sqsubseteq$ . We want to prove safety properties over, say **rm**, so any  $\hat{\mathbf{rm}}$  that concretizes to a set of states that contains the “computed”  $\hat{\mathbf{rm}}' = \hat{\sigma}'_- \hat{\neq}^{rm} \hat{\sigma}''_+$ , is sound. This has the practical effect of allowing weakening within the abstract interpretation.

Now, we turn our attention to discussing  $\hat{\neq}$ . The idea is that this abstracts removed or added states. Recall the definition of removed states from our concrete semantics,  $\mathbf{rm} = \{\sigma'_- \in \Sigma'_- \mid \sigma'_- \neq^{rm} \Sigma'_+\}$ , and unrolling  $\neq$  (definition 1) gives us  $\mathbf{rm} := \{\sigma_- \in \Sigma_- \mid \forall \sigma_+ \in \Sigma_+. \sigma_- \neq \sigma_+\}$ .

The soundness of  $\hat{\neq}$  is defined as follows.

$$\{\sigma \in \gamma(\hat{\sigma}) \mid \sigma \neq \gamma(\hat{\sigma}')\} \subseteq \gamma(\sigma \hat{\neq} \hat{\sigma}')$$

The glaring problem with such an abstraction is that, usually, there is no way to “iterate” through abstract elements. We can naturally lift  $\sim$  to our abstract domain,  $\hat{\sim}$  (and in turn  $\hat{\neq}$ ), by defining it on our abstract elements instead of concrete states. However, if we have no way to “iterate” through  $\hat{\sigma}$ , then we must characterize an *entire*  $\hat{\sigma}_-$  as either removed or not removed. This simple abstraction of removed states defines  $\hat{\sigma} \hat{\neq} \hat{\sigma}'$  as  $\hat{\sigma} \hat{\neq} \hat{\sigma}'$ . We can see this as the concrete definition of **rm** with only one element in  $\Sigma_-$  and  $\Sigma_+$ .

This is quite a coarse abstraction of removed and added states (all or none). Furthermore, we notice that our **rm** property is interpreted over old states, but requires the entire set of new states to compute it. However, this simple **rm** abstraction actually reduces **rm** properties to no longer require the set of new states. Instead, it can be formulated as a property dependent only on state pairs. Since with this abstraction either all states will be in **rm** or none, it is equivalent to checking  $\sigma \sim \sigma'$  or  $\sigma \not\sim \sigma'$  over reachable pairs. These types of properties can be proven in many existing frameworks, so we note that the additional properties required of  $\mathcal{D}$  are essential to our approach.

Now, we present two ways to implement  $\hat{\neq}$ . Option 1 places a larger restriction on  $\mathcal{D}$ , while option 2 is more permissive but most natural for symbolic domains. While 2 could probably be implemented for most domains, it will likely require an additional abstraction to some symbolic domain. The instantiation of our parameterized abstract interpreter presented in section 4 will use option 2 with a symbolic domain.

- (1) We can achieve a more precise  $\hat{\neq}$ , that cannot be reduced to a predicate over pairs of states, if  $\mathcal{D}$  supports finite partitioning. A finitely partitioned abstract domain is such that  $\hat{\sigma} = \hat{\sigma}_1, \hat{\sigma}_2, \dots, \hat{\sigma}_n$  where  $\gamma(\hat{\sigma}_i) = \Sigma_i$  and  $\gamma(\hat{\sigma}) = \gamma(\hat{\sigma}_1) \cup \gamma(\hat{\sigma}_2) \cup \dots \cup \gamma(\hat{\sigma}_n)$ . Then, we may

491 define  $\hat{\sigma} \hat{\approx} \hat{\sigma}'$  as  $\forall \hat{\sigma}_i \in \hat{\sigma}. (\forall \hat{\sigma}'_i \in \hat{\sigma}'. \hat{\sigma}_i \hat{\neq} \hat{\sigma}'_i)$  where  $\hat{\neq}$  is the  $\neq$  relation lifted to our abstract  
 492 domain (defining the operator on the new type). Such a domain exists.

- 493 (2) For domains which are not finitely partitioned, hope is not lost. Let's consider abstracting  
 494 removed states with  $\hat{\approx}^{\text{rm}}$ . If we had an abstract element, say  $\hat{\sigma}'_-$  that represented all possible  
 495 states where for any one of those states, say  $\sigma'_- \in \gamma(\hat{\sigma}'_-)$ ,  $\forall \sigma_+ \in \gamma(\hat{\sigma}_+).$   $\sigma'_- \not\sim^{\text{rm}} \sigma_+$  held;  
 496 then we could meet that with the reachable state element to achieve an abstraction of **rm**,  
 497  $\hat{\sigma}_- \sqcap \hat{\sigma}'_-$ . This meet would represent all reachable states that also satisfy our set lifted state  
 498 correspondence relation (*i.e.* removed condition).

499 Now, we are left with constructing a  $\hat{\sigma}'_-$  such that  $\forall \sigma'_- \in \gamma(\hat{\sigma}'_-).$   $(\forall \sigma_+ \in \gamma(\hat{\sigma}_+).$   $\sigma'_- \not\sim^{\text{rm}} \sigma_+)$ . First, we make the following observation.  $\hat{\sigma}$  is a state predicate, we even defined its  
 500 concretization in terms of all concrete states which model it. We can lift this state predicate,  
 501 say  $\hat{\sigma}_+$ , to a global property by closing it with respect to variables modified by  $s_+$ . If we,  
 502 say, universally quantify over these modified variables, then we achieve a property (that is  
 503 not a state predicate) which holds for all possible valuations of the modified variable, *i.e.* a  
 504 property that represents all states which  $\hat{\sigma}_+$  holds.

505 A representation of all pairs of states which satisfies  $\neq$ , say  $\hat{\sigma}_{sc}$ , is a predicate on pairs of  
 506 states. We could do the same trick to remove the dependency on one of the states, lowering  
 507 it to a single state predicate. Now, we have the tools we need to construct our  $\hat{\sigma}'_-$ . We let  $\eta_+$   
 508 be a function which given a predicate on *at least*  $\sigma_+$ , say  $\varphi$ , removes the dependency on old  
 509 states by producing a new assertion that holds for all  $\sigma_+$  which  $\varphi$  represented.

510 DEFINITION 12. *An abstraction of the removed state operator ( $\hat{\approx}^{\text{rm}}$ ) for domains that are not  
 511 finitely partitioned. We require  $\hat{\sigma}_{sc}$  and  $\eta_+$  such that:*

$$\begin{aligned} 512 \{(\sigma_-, \sigma_+) \mid \sigma_- \not\sim^{\text{rm}} \sigma_+\} &\subseteq \gamma(\hat{\sigma}_{sc}) \\ 513 \{\sigma_- \mid \forall \sigma_+ \in \gamma(\hat{\sigma}_+). (\sigma_-, \sigma_+) \models \hat{\sigma}_{sc}\} &\subseteq \gamma(\eta_+(\hat{\sigma}_+ \sqsubseteq \hat{\sigma}_{sc})) \end{aligned}$$

514  $\eta_+(\hat{\sigma}_+ \sqsubseteq \hat{\sigma}_{sc})$  can equivalently be defined as all  $\sigma_-$  such that  $\gamma(\hat{\sigma}_+) \subseteq \gamma(\hat{\sigma}_{sc})$  and  $\sigma_- \in$   
 515  $\{\sigma'_- \mid (\sigma'_+, \sigma'_-) \in \gamma(\hat{\sigma}_{sc})\}$  holds.

516 In order to define  $\hat{\sigma}_{sc}$ , a definition  $\hat{\sim}$  will be required. This will be used to compute  $\hat{\text{rel}}$ . We  
 517 require the same constructions for  $\hat{\approx}^{\text{add}}$ .

518 Now, we have  $\hat{\sigma}_- \hat{\approx}^{\text{rm}} \hat{\sigma}_+$  if and only if  $\hat{\sigma}_- \sqcap \eta_+(\hat{\sigma}_+ \sqsubseteq \hat{\sigma}_{sc})$

519 DEFINITION 13. *An instantiation of  $\eta$  from definition 12 for symbolic domains which include  
 520 universal quantification in the assertion language.*

$$521 \eta_+(\hat{\sigma}_2) := \forall \text{mod}(s_+). \hat{\sigma}_2$$

522 Where  $\text{mod}(s_+)$  is all variables modified by  $s_+$ .

## 523 4 Instantiation with Variable Correspondence, Hoare Logic, Incorrectness Logic, and 524 IMP

525 Now, we will instantiate the parameterized abstract interpreter given in section 3.5.

526 Our underlying programming language will be IMP with non-deterministic assignment whose  
 527 syntax is given in section A. The concrete semantics of imp is given with the judgment form  
 528  $\sigma \vdash s \Downarrow^* \sigma'$ .  $\sigma$  is a map from program variables to concrete values, which are either numbers or  
 529 booleans.

530 Now, we will define an instantiation of state correspondence for IMP, which we will call the  
 531 variable correspondence. The grammar and semantics are given below.

540     DEFINITION 14. An instantiation of state correspondence for IMP, which we will refer to as the  
 541     variable correspondence.

$$\begin{array}{ll} \text{variable correspondence} & vc ::= (u), (u') \\ \text{universional state bit} & u ::= u, x \mid x' \\ \text{variables} & x \end{array}$$

546     With the variable correspondence,  $\sim^{rm} = \sim^{add}$ . Furthermore, a vc, say  $u, u'$  is well formed only if  $u$   
 547     has the same size as  $u'$ . However,  $u$  need not have the same number of unique variables as  $u'$ .

548     For a given  $vc := (x_1, \dots, x_n), (x'_1, \dots, x'_n)$ , we define  $\sim$  as follows.

$$\sigma_- \sim \sigma_+ \text{ if and only if } (\sigma_-\langle x_1 \rangle, \dots, \sigma_-\langle x_n \rangle) \neq (\sigma_+\langle x'_1 \rangle, \dots, \sigma_+\langle x'_n \rangle)$$

551     Where  $\sigma\langle x \rangle$  looks up the variable  $x$  in state  $\sigma$ .

552     Now, we will instantiate our over approximating abstract semantics with Hoare logic.

$$\hat{\sigma} \vdash s \Downarrow \hat{\sigma}' := \vdash \{\hat{\sigma}\} s \{ \hat{\sigma}' \}$$

554     And our under approximating abstract semantics with Incorrectness logic.

$$\hat{\sigma} \vdash s \Downarrow \hat{\sigma}' := \vdash [ \hat{\sigma} ] s [ \hat{\sigma}' ]$$

556      $\hat{\sigma}$  will be formulae in an assertion language which is first order logic along with IMP expression.  
 557      $\hat{\sigma}_-$  will include variables subscripted with  $-$ , and likewise for  $\hat{\sigma}_+$  and  $+$ . Abstract elements which  
 558     represent state pair predicates may contain variables with both  $-$  and  $+$  subscripts.

559     Now, all that is left is to define is  $\hat{\sigma}$ . First, for the  $vc := (x^1, \dots, x^n), (y^1, \dots, y^n)$  we define  
 560      $\hat{\sigma}_{sc}$  as  $(x^1, \dots, x^n) \neq (y^1, \dots, y^n)$ , and  $\hat{\sigma}'_{sc}$  as  $(x^1, \dots, x^n) = (y^1, \dots, y^n)$ . We define  $\hat{\sigma} \hat{\neq} \hat{\sigma}'$  as  $\hat{\sigma} \wedge$   
 561      $(\forall mod(s). \hat{\sigma}' \implies \hat{\sigma}_{sc})$  when  $\hat{\sigma}'$  represents reachable states of program  $s$ . This leaves us with an  
 562     instantiated abstract interpreter.

563     DEFINITION 15. Our abstract interpreter with IMP, vc, Hoare logic, and Incorrectness logic.

$$\begin{array}{l} \vdash \{\hat{\sigma}_-\} s_- \{\hat{\sigma}'_-\} \vdash \{\hat{\sigma}_+\} s_+ \{\hat{\sigma}'_+\} \vdash [\hat{\sigma}_-] s_- [\hat{\sigma}''_-] \vdash [\hat{\sigma}_+] s_+ [\hat{\sigma}''_+] (\hat{\sigma}'_- \wedge \hat{\sigma}'_+ \wedge \hat{\sigma}'_{sc}) \implies \hat{rel} \\ \vdash [\hat{\sigma}'_- \sqcap (\forall mod(s_+). \hat{\sigma}''_+ \implies \hat{\sigma}_{sc})] \implies \hat{rm} \quad [\hat{\sigma}'_+ \sqcap (\forall mod(s_-). \hat{\sigma}''_- \implies \hat{\sigma}_{sc})] \implies \hat{add} \end{array}$$


---


$$\hat{\sigma}_-, \hat{\sigma}_+ \vdash s_- \pm s_+ \Downarrow^\# \hat{rm}, \hat{add}, \hat{rel}$$

569     The major difference from our parameterized abstract interpreter is the computation of  $\hat{rel}$ .  
 570     We had  $[(\hat{\sigma}'_- \sim^{rm} \hat{\sigma}'_+) \sqcup (\hat{\sigma}'_- \sim^{add} \hat{\sigma}'_+)] \sqsubseteq \hat{rel}$ . Since  $\sim^{rm} = \sim^{add}$  and  $\sqsubseteq$  is  $\implies$  in our domain,  
 571     we first simplify to  $(\hat{\sigma}'_- \hat{\sim} \hat{\sigma}'_+) \implies \hat{rel}$ . Now, the soundness of  $\hat{\sim}$  is defined as  $\{(\sigma'_-, \sigma'_+) \in$   
 572      $\gamma(\hat{\sigma}'_-) \times \gamma(\hat{\sigma}'_+) \mid \sigma'_- \sim \sigma'_+\} \subseteq \gamma(\hat{\sigma}'_- \hat{\sim} \hat{\sigma}'_+)$ . We notice that  $\gamma(\hat{\sigma}'_- \wedge \hat{\sigma}'_+) = \{\gamma(\hat{\sigma}'_-) \times \gamma(\hat{\sigma}'_+)\}$  and  
 573      $\gamma(\hat{\sigma}'_{sc}) = \{(\sigma_-, \sigma_+) \mid \sigma_- \sim \sigma_+\}$ . Thus,  $\gamma(\hat{\sigma}'_- \wedge \hat{\sigma}'_+ \wedge \hat{\sigma}'_{sc}) = \{(\sigma'_-, \sigma'_+) \in \gamma(\hat{\sigma}'_-) \times \gamma(\hat{\sigma}'_+) \mid \sigma'_- \sim \sigma'_+\}$ .

## 5 Evaluation

576     In this section, we use the abstract interpreter defined in definition 15 to prove properties about  
 577     IMP program revisions. We will use very simple examples that illustrate the coherence of our  
 578     constructions with respect to program difference properties.

580     First, let's consider  $x := 1 \pm y := 2$ . With  $vc := (x), (y)$ , we see that all old states ( $x_- = 1$ ) are  
 581     removed, and all new states ( $y_+ = 2$ ) are added. No reachable states satisfy the state correspondence  
 582     ( $x_- = y_+$ ), so  $\hat{rel}$  is empty. This is all quite intuitive, but what about about  $vc := (x), (x)$ ? Let's  
 583     check if  $\hat{add} := (y_+ = 2)$  is derivable. We have  $\vdash [\top] x := 1 [x_- = 1]$  and  $\vdash \{\top\} y := 2 \{y_+ = 2\}$ .  
 584     We generate the verification condition:  $[y_+ = 2 \wedge (\forall x_- x_- = 1 \implies x_- \neq x_+)] \implies y_+ = 2$ . This is  
 585     certainly satisfiable, so  $y_+ = 2$  is a sound over approximation of  $\hat{add}$ . However, the set of reachable  
 586     new states will always be a sound over approximation of  $\hat{add}$ , which we can clearly see with the  
 587     formula above.

589 Let's see if we can pin down **add** any better. Well, the antecedent of our verification condition  
 590 precisely describes **add**. It reduces to  $y_+ = 2 \wedge x_+ \neq 1$ . This makes sense, as it will be true of all  
 591 added states.

## 594 6 Related Work

595 Verifying properties of the difference between two programs is not new. Relational Hoare Logic  
 596 (RHL) [3] defines a means to reason about the relationship between two programs. It is used to  
 597 prove many categories of properties which frame and formalize the semantics of a code change  
 598 such as regression verification [4], translation validation [7], contextual equivalence [8], refinement,  
 599 and others [2]. The strategy used to prove these properties is often either the construction of a  
 600 product program [1] to reduce the problem to a functional safety property, or to rely on similarities  
 601 between the two programs.

602 Some work has been done to develop analysis tools and techniques particularly for program  
 603 diffs and their properties of interest. Semantic Diff [5] and differential symbolic execution [10] are  
 604 tools which summarize the behavioral differences between two version of a program. Differential  
 605 assertion checking [6] checks if assertion which holds in one program does or does not hold in  
 606 another. There is also work [re-find citation] using sound abstract interpretation techniques to  
 607 prove output equivalence given the same inputs.

608 The related work which most closely resembles ours is “Abstract Semantic Differencing for  
 609 Numerical Programs” [9]. It presents an abstract interpretation framework that over approximates  
 610 program differences. This final abstract state of this tool is a list of abstracted guard conditions  
 611 and the abstract values of correlated variables under certain guard conditions. This can be seen as  
 612 summarizing the differences, providing bounds for those differences, and providing conditions for  
 613 the program variables to be equal. This work also has the concept of added, removed, and related  
 614 states; however, these terms are used to mean a partition on the final abstract state. So, for example,  
 615 a new value for a certain variable will be in the added state. This work also relies on a product  
 616 program, variable correspondence, and correlating abstract domain.

## 619 7 Conclusion

620 We introduced a concrete semantics that describes the semantic effect of program revisions which  
 621 characterizes reachable states as removed, related, or added. These sets are created via a parameter-  
 622 ized equivalence relation on states which ensures the semantics are grounded in the intuitive notion  
 623 developers have about program differences. We saw the intricacies associated with proving (safety)  
 624 program difference properties which necessitated precise abstract domains, and both over and  
 625 under approximating abstract semantics. A heavily parameterized abstract interpreter for program  
 626 revision properties was introduced, then an instantiation of it with Hoare Logic and Incorrectness  
 627 Logic was used to prove some simple examples.

## 630 Acknowledgments

631 This report is based on a project advised by Bor-Yuh Evan Chang with help from Matthew Hammer.  
 632 Of course, this project would not exist without them. The idea originated from a grant submitted by  
 633 Evan Chang, and their insights, guidance, and contributions have been invaluable in the creation  
 634 of this paper. This paper is not solely authored by me, but since it is a final project, I will instead  
 635 acknowledge that here.

## 638 A IMP

639 DEFINITION 16. *IMP grammar.*

```

640     expressions   e ::= n | x | e1 + e2 | e1 - e2 | e1 ÷ e2 | e1 % e2
641     boolean expressions   b ::= bv | e1 = e2 | e1 != e2 | b1 && b2 | b1 || b2 | e1 >= e2 | e1 < e1 |
642         values   v ::= n
643         commands  s ::= x := e | if b then s1 else s2 |
644                           s1 ; s2 | skip | while(b){s} | x := !
645         variables  x
646         numbers   n
647         booleans  bv
648

```

## 649 References

- 650 [1] Gilles Barthe, Juan Manuel Crespo, and César Kunz. 2011. Relational Verification Using Product Programs. In *FM 2011: Formal Methods*, Michael Butler and Wolfram Schulte (Eds.). Springer, Berlin, Heidelberg, 200–214. [doi:10.1007/978-3-642-21437-0\\_17](https://doi.org/10.1007/978-3-642-21437-0_17)
- 653 [2] Bernhard Beckert and Mattias Ulbrich. 2018. Trends in Relational Program Verification. In *Principled Software Development: Essays Dedicated to Arnd Poetzsch-Heffter on the Occasion of his 60th Birthday*, Peter Müller and Ina Schaefer (Eds.). Springer International Publishing, Cham, 41–58. [doi:10.1007/978-3-319-98047-8\\_3](https://doi.org/10.1007/978-3-319-98047-8_3)
- 655 [3] Nick Benton. [n. d.]. Simple Relational Correctness Proofs for Static Analyses and Program Transformations. ([n. d.]).
- 656 [4] Benny Godlin and Ofer Strichman. [n. d.]. Regression Verification: Proving the Equivalence of Similar Programs. ([n. d.]).
- 658 [5] Jackson and Ladd. 1994. Semantic Diff: a tool for summarizing the effects of modifications. In *Proceedings 1994 International Conference on Software Maintenance*. 243–252. [doi:10.1109/ICSM.1994.336770](https://doi.org/10.1109/ICSM.1994.336770)
- 660 [6] Shuvendu K. Lahiri, Kenneth L. McMillan, Rahul Sharma, and Chris Hawblitzel. 2013. Differential assertion checking. In *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering*. ACM, Saint Petersburg Russia, 345–355. [doi:10.1145/2491411.2491452](https://doi.org/10.1145/2491411.2491452)
- 662 [7] Nuno P. Lopes, David Menendez, Santosh Nagarakatte, and John Regehr. 2015. Provably correct peephole optimizations with alive. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation*. ACM, Portland OR USA, 22–32. [doi:10.1145/2737924.2737965](https://doi.org/10.1145/2737924.2737965)
- 664 [8] Andrzej S Murawski, Steven J Ramsay, and Nikos Tzevelekos. [n. d.]. Game Semantic Analysis of Equivalence in IMJ. ([n. d.]).
- 666 [9] Nimrod Partush and Eran Yahav. 2013. Abstract Semantic Differencing for Numerical Programs. In *Static Analysis*, David Hutchison, Takeo Kanade, Josef Kittler, Jon M. Kleinberg, Friedemann Mattern, John C. Mitchell, Moni Naor, Oscar Nierstrasz, C. Pandu Rangan, Bernhard Steffen, Madhu Sudan, Demetri Terzopoulos, Doug Tygar, Moshe Y. Vardi, Gerhard Weikum, Francesco Logozzo, and Manuel Fähndrich (Eds.). Vol. 7935. Springer Berlin Heidelberg, Berlin, Heidelberg, 238–258. [doi:10.1007/978-3-642-38856-9\\_14](https://doi.org/10.1007/978-3-642-38856-9_14) Series Title: Lecture Notes in Computer Science.
- 670 [10] Suzette Person, Matthew B. Dwyer, Sebastian Elbaum, and Corina S. Păsăreanu. 2008. Differential symbolic execution. In *Proceedings of the 16th ACM SIGSOFT International Symposium on Foundations of software engineering (SIGSOFT '08/FSE-16)*. Association for Computing Machinery, New York, NY, USA, 226–237. [doi:10.1145/1453101.1453131](https://doi.org/10.1145/1453101.1453131)