

Abstract Semantic Differencing for Numerical Programs

Abstract

We address the problem of computing semantic differences between a program and a patched version of the program. Our goal is to obtain a precise characterization of the difference between program versions, or establish their equivalence when no difference exists.

We focus on computing semantic differences in numerical programs where the values of variables have no a-priori bounds, and use abstract interpretation to compute an over-approximation of program differences. Computing differences and establishing equivalence under abstraction requires abstracting relationships between variables in the two programs. Towards that end, we first construct a *correlating program* in which these relationships can be tracked, and then use a *correlating abstract domain* to compute a sound approximation of these relationships. To better establish equivalence between correlated variables and precisely capture differences, our domain has to represent non-convex information. To balance precision and cost of this representation, our domain may over-approximate numerical information as long as equivalence between correlated variables is preserved.

We have implemented our approach in a tool called DIZY, built on the LLVM compiler infrastructure and the APRON numerical abstract domain library, and applied it to a number of challenging real-world examples, including programs from the GNU core utilities, Mozilla Firefox and the Linux Kernel. Our evaluation shows that DIZY often manages to establish equivalence, describes precise approximation of semantic differences when difference exists, and reports only a few false differences.

1. Introduction

Understanding the semantic difference between two versions of a program is invaluable in the process of software development. A developer making changes to a program is often interested in answering questions like: (i) did the patch add/remove the desired functionality? (ii) does the patch introduce other, *unexpected*, behaviors? (iii) which regression tests should be run? Answering these questions manually is difficult and time consuming.

Semantic differencing has received much attention in classical work (e.g., [7, 8]) and has recently seen growing interest for various applications ranging from testing of concurrent programs [3], understanding software upgrades [10], to automatic generation of security exploits [2].

Existing Techniques Existing techniques mostly offer under-approximating solutions, the prominent of which is regression testing which provides limited assurance of behavior equivalence. Furthermore, running regression on large systems requires significant time and compute resources. Godlin and Strichman [6] rely on bounded model checking techniques to produce a (binary) result regarding (input-output) equivalence of two closely related numerical programs. Kawaguchi et al. [11] define the notion of **conditional equivalence but cannot compute differences in loops**. Other techniques for computing semantics differences [13, 14] rely

on symbolic execution, may miss differences, and generally unable to prove equivalence.

We present an approach based on abstract program interpretation [4] for a *sound*, succinct representation of changed program behaviors and proving equivalence. Our method focuses on abstracting relationships between variables, and therefore behaviors, in both versions allowing us to achieve a precise description of difference and prove equivalence while ignoring other program information which may encumber a traditional analysis but is less relevant in our setting¹.

Problem Definition We define the problem of *semantics differencing* as follows: Given a pair of programs (P, P') which agree on the number and type of inputs, for every execution π of P that originate from an input i and a corresponding execution π' of P' that originates from the *same input* i our goal is:

- Check whether π and π' agree on output i.e., are output-equivalent.
- In case of difference in behavior, provide a description of difference.

We intentionally define the notion of input and output equivalence loosely at this point, and we discuss several realizations of these in later sections.

To answer the question of semantic differencing for infinite-state programs, we employ abstract interpretation. Though the notion of difference is well defined in the concrete case, defining and soundly computing it under abstraction is challenging as:

- Differencing requires correlation of *different program executions* meaning the abstraction must be able to capture input-equivalent executions, and distinguish ones that are not input-equivalent.
- Establishing equivalence of abstract output values does not entail equality between the concrete output value they represent.

To address these challenges, we introduce two new concepts: (i) a *correlating program*, a single program $P \bowtie P'$ that captures the behaviors of both P and P' in a way that facilitates abstract interpretation; (ii) a *correlating abstract domain*, tracking relationships between variables in P and variables in P' by tracking equivalences in $P \bowtie P'$.

Correlating Program Abstracting relationships allows us to maintain focus on difference while omitting (whenever necessary for scalability) parts of the behavior that does not entail difference. In order to monitor these relationships we created a *correlating program* which captures the behavior of both the original program and its patched version. Instead of designing a correlating semantics that is capable of co-executing two programs, we chose to automatically construct the correlating program such that we can benefit from the use of standard analysis frameworks for analyzing the

¹ Eran: you didn't like this sentence, i love it, i'm trying to say we abstract away numerical information etc. and focus on relationships, thus we have better results since we are not encumbered by this less relevant information

```

int sign(int x) {
  int sgn;
  if (x < 0)
    sgn = -1
  else
    sgn = 1
  return sgn
}

int sign'(int x) {
  int sgn;
  if (x < 0)
    sgn = -1
  else
    sgn = 1
  if (x==0)
    sgn = 0
  return sgn
}

```

Figure 1: Two simple implementations of the *sign* operation.

resulting program. Another advantage of this new construct, is that you may apply other methods for equivalence checking directly on it [14] as the correlation allows for a much more fine-grained equivalence checking (between local variables and not only output).

Correlating Abstraction Our abstraction holds data of both sets of variables, joined together and is initialized to hold equality over all matched variables. This means we can reflect relationships without necessarily knowing the actual value of a variables (we can know that $x_{old} = x_{new}$ even though actual values are unknown). We ran our analysis over the correlating program while updated the domain to reflect program behavior. Since some updates may result in non-convex information (e.g. taking a condition of the form $x \neq 0$ into account), our domain has to represent non-convex information, at least temporarily. We address this by working with a powerset domain of a convex representation with partitioning according to equivalence criteria to avoid exponential blowup. Our domain may over-approximate numerical information as long as equivalence between correlated variables is preserved.

1.1 Main Contributions

The main contributions of this paper are as follows:

- we present a method for abstract interpretation of a pair of programs (P, P') for *sound* semantic equivalence and differencing by abstracting direct relationships between (P, P') variables in a partially disjunctive domain. We describe a partitioning technique for state reduction and scaling. We define a widening operator for abstracting unbound paths in our domain.
- we phrase a new technique for syntactically interleaving a pair of programs (P, P') for the creation of a *correlating program* $P \bowtie P'$ which contains the semantics of both programs. We propose an analysis over the program for characterizing program equivalence and difference, based on the aforementioned abstraction, given the properties of the correlating program which aligns (P, P') executions.
- We have implemented our approach in a tool based on the LLVM compiler infrastructure and the APRON numerical abstract domain library, and evaluated it using select patches from open-source software including GNU core utilities, Mozilla Firefox, and the Linux Kernel. Our evaluation shows that the tool often manages to establish equivalence, reports useful approximation of semantic differences when differences exists, and reports only a few false differences.

2. Overview

Consider the simple example program of Fig. 1, inspired by an example from [15]. For this example, we would like to establish that the output of *sign* and *sign'* only differ in the case where $x = 0$ and that the difference is $sgn = 1 \neq sgn' = 0$. An optimal *abstract* characterization of behavior is as following:

$$\begin{aligned}
\sigma^1 &= \{x < 0, sgn \mapsto -1\} & \equiv & \sigma'^1 = \{x' < 0, sgn' \mapsto -1\} \\
\sigma^2 &= \{x = 0, sgn \mapsto 1\} & \neq & \sigma'^2 = \{x' = 0, sgn' \mapsto 1\} \\
\sigma^3 &= \{x > 0, sgn \mapsto 1\} & \equiv & \sigma'^3 = \{x' > 0, sgn' \mapsto 1\}
\end{aligned}$$

The abstraction precisely captures and describes equivalence between σ^1, σ^3 and their tagged counterpart for $x < > 0$ while the difference is captured by $\sigma^2 \neq \sigma'^2$ when $x = 0$.

As a first naive attempt to achieve this description, one could try to analyze each version of the program separately and compare the (abstract) results. However, this is clearly unsound, as equivalence under abstraction does not entail concrete equivalence. For example, using an interval analysis [?] would yield that in both programs the value of *sgn* ranges in the same interval $[-1, 1]$, missing the fact that *sign* never returns the value 0 as depicted in Fig. ??.

Furthermore, this result entirely ignores how x affects the value of *sgn* thus we would have no means to differentiate correctly, according to input (e.g. we will get the same result for the $-1 * sign$ function).

To establish equivalence under abstraction, we need to abstract relationships between the values of variables in *sign* and *sign'* under the assumption of equivalence of input. Specifically, we need to track the relationship between the values of *sgn* in both versions and see whether we can establish their equivalence. Tracking relationships dictates performing a *joint analysis* that employs a *correlating abstraction* allowing us to bind variables of both programs in one abstract state.

A correlating-oriented abstraction is well suited for proving equivalence as it allows focusing on relationships between versions of variables while abstracting away other (numerical) information allowing us to scale better. Most importantly, such an abstraction guarantees that equivalence will be reported soundly: as in a separate analysis we abstracted $\langle sgn \mapsto -1 \rangle$ and $\langle sgn \mapsto 1 \rangle$ towards an interval $\langle sgn \mapsto [-1, 1] \rangle$, and again for *sgn'* values, separately, which cannot assure equivalence. Instead we abstract $\langle sgn = sgn' \mapsto -1 \rangle$ and $\langle sgn = sgn' \mapsto 1 \rangle$ as $\langle sgn = sgn' \rangle$ which soundly assures equivalence although all other variable information has been abstracted away.

We present an abstraction over dual program state, thats able to correlate paths that originate from the same input as well as produce a characterization of equivalence and difference which reflects change in behavior precisely. For example, we produce the following constraints for *sign* and *sign'*:

$$\begin{aligned}
\sigma_x^1 &= \{x = x' < 0, sgn = sgn' \mapsto -1\} \\
\sigma_x^2 &= \{x = x' = 0, sgn \mapsto 1, sgn' \mapsto -1\} \\
\sigma_x^3 &= \{x = x' > 0, sgn = sgn' \mapsto 1\}
\end{aligned}$$

To arrive at this result, we initially abstracted the dual program state by analyzing both programs sequentially $(P; P')$, while updating the shared state with data regarding both sets of variables. We allow direct relationships between versions of variables, this will be of utmost importance later on when we over approximate paths for scalability as it will allow us to maintain equivalence soundly. In order to correlate paths by input and arrive at a precise disjunction, the analysis initially assumes input equivalence $\vec{i} = \vec{i}'$.

As we advance through the analysis of P , we will accumulate the disjunction of all possible path constraints in its final state (this is similar to trace partitioning [15]). At this point, as we continue to analyze P' , each disjunct representing a path in P will be further conjuncted with all of P' paths. This will produce a precise disjunction for differencing as each path in P will be split and conjuncted with all of P' paths, while avoiding considering conjunctions that disagree on input due to our input equivalence assumption. An illustration of the joint analysis for the *sign* example can be seen in Fig. 2 including markings for feasible and infeasible paths.

Essentially, our analysis aims to establish correspondence between paths in P and P' by first analyzing all of P paths and then

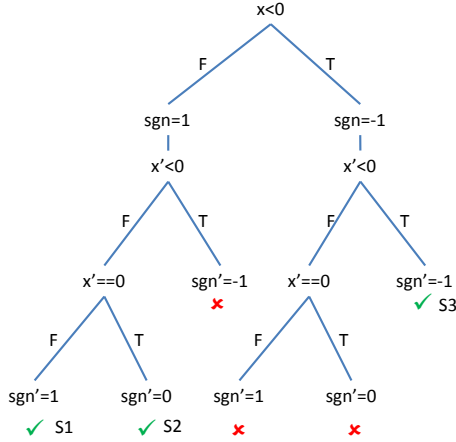


Figure 2: Joint $sign; sign'$ analysis

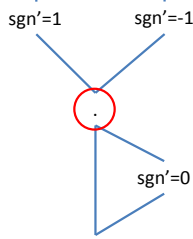


Figure 3: $sign \boxtimes sign'$ analysis

attempting to correlate with P' paths. Clearly, this abstraction is unfeasible for most cases as the number of paths to be considered is exponential (we defer the case of loops where this number is unbound). Therefore we refine our abstraction by using a partially disjunctive domain, partitioned by *equivalence criteria*.

As the goal of work is to distinguish equivalent from differencing behaviors, using equivalence as criteria for merging paths is apt. The partitioning will abstract together paths that hold equivalence for the same set of variables, allowing for a maximum of $2^{|V|}$ disjunctions in the abstract state, where V is the set of correlated (output?) variables. This criteria can be refined, by adding to V , to provide a more precise result or alternatively can become more coarse by allowing only certain equivalence classes of 2^V .

For example partitioning the result of Fig. 2 according to our criteria would abstract behaviors $S1$ and $S3$ together, as they hold equivalence for sgn . The merge would abstract away data regarding x and represent sgn as the $[-1, 1]$ interval, losing precision but gaining reduction in state size. This lose of precision is acceptable as it is complemented by the offending state $S2$. Still, not much is gained from this partitioning, as it is performed at the final state, where we may have already reached an exponential amount of disjunctions.

To truly gain a reduction of state size, we must perform partitioning dynamically, as the analysis is executed i.e. at earlier program locations. This cannot be achieved using a sequential composition $P; P'$. Looking at Fig. 2 we immediately see that equivalence holds only at final states. Intuitively, this is caused due to a command in P having to "wait" for it's equivalent command to arrive in P' . To overcome this, we present the correlating program $P \boxtimes P'$

```
int sign(int x) {
  int x' = x;
  guard g1 = (x < 0);
  guard g1' = (x' < 0);
  int sgn;
  int sgn';
  if (g1) sgn = -1;
  if (g1') sgn' = -1;
  if (!g1) sgn = 1;
  if (!g1') sgn' = 1;
  guard g2' = (x' == 0);
  if (g2') sgn' = 0;
}
```

Figure 4: Correlating program $sign \boxtimes sign'$.

```
int sum(int arr[], unsigned len) {
  int result = 0;
  for (unsigned i = 0; i < len; i++)
    result += arr[i];
  return result;
}
```

Figure 5: A simple looping program for array summation.

which allows for earlier partitioning by "saving the need to wait" as it interleaves P and P' commands in an optimized manner, and informs the analysis that it need not wait any further and partitioning is permitted. Fig. 3 depicts the analysis of $sign \boxtimes sign'$ (shown in Fig. 4) where the partitioning location is marked in red. We define these partitioning locations as *correlation points* (denoted CP) and they are a sub product of the correlating program build process. We will further describe the specifics of creating $P \boxtimes P'$ in Section 6 and only shortly say that the interleaving is chosen according to a syntactic diff process over a guarded command language version of the programs.

Although we achieved a reduction in state size using partitioning, we have yet to account for programs with an unbound number of paths, created by loops. Unbound path lengths means a potentially unbound analysis as all paths are abstracted. This is mainly where previous approaches fall short [1]. To overcome this, we define a widening operator for our domain, based on the convex sub-domain widening operator. The main challenge here, as our state is a set of convex objects, is finding an optimal pairwise matching between objects for a precise widened result. Optimally, we would like to pair objects that adhere to the same "looping path" meaning we would want to match π_i 's abstraction with a π_{i+1} that results from taking another step in the loop. This basically requires encoding path information along with the sub-state abstraction. This information is acquired by simply keeping guard values explicitly, as they appear in our correlating program, inside the state. As guard values (true or false) reflect branch outcomes, they can be used to match sub-states that advanced on the loop by matching their guard values (for easier matching and better precision we separate guards from other variables in our implementation).

We note that the correlating program is significant to maintaining equivalence over loops. To demonstrate this we perform the simple exercise of checking equivalence of a small looping program with itself. Consider the array summation program in Fig. 5. The state of the art in program equivalence cannot prove that this small program is equivalent to itself. To emphasize the importance of the correlating program, we will first show the result of an analysis of $sum; sum$ which will be:

$$\begin{aligned} \sigma_x^1 &= \{len = len' \mapsto 0, result = result' \mapsto 0\} \\ \sigma_x^2 &= \{len = len' > 0, i \leq len, i' \leq len'\} \end{aligned}$$

```

int sum(int arr[], unsigned len) {
    unsigned len' = len;
    int arr'[] = arr;
    int result = 0;
    int result' = 0;
    {
        unsigned i = 0;
        unsigned i' = 0;
    l: guard g = (i < len);
    l': guard g' = (i' < len');
        if (g) result += arr[i];
        if (g') result' += arr'[i'];
        if (g) i++;
        if (g') i'++;
        if (g) goto l;
        if (g') goto l';
    }
}

```

Figure 6: $sum \bowtie sum$

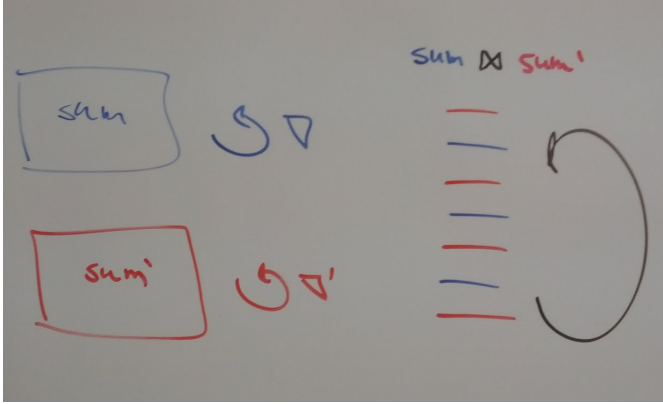


Figure 7: Widening $sum; sum$ vs. $sum \bowtie sum$

This loss of equivalence occurred due to the inability precisely track the relationship of $result$ and $result'$ over $sum; sum$. As we widened the first loop to converge, all paths passing through that loop were merged together, losing the ability to be "matched" with the second loop waiting further down the road. Performing the same analysis on $sum \bowtie sum$ instead (seen in Fig. ??, allows the maintaining of equivalence, as the loops are interleaved correctly to allow establishing $result = result'$ as a loop invariant, surviving the widening process to prove equivalence at the end as the result would be:

$$\sigma_{\bowtie}^1 = \{len = len' \geq 0, result = result'\}$$

The conceptual difference of these two analyses is depicted in Fig. 7.

Definition of difference So far, we defined difference in programs as difference in output variable values at the program final state. Our work extends the notion of difference beyond that, allowing for several program locations to be identified as viable for differentiation (we name these locations *differencing points* denoted *DP*). Other than the programs endpoints, we include:

- Emitting an output value (using a system function or through a global variable)
- Array access

This means that procedures may now differ, even if they hold equivalence on the final state. Thus we define difference also as producing a different output value in mid-execution or accessing an array

through a different index. This means that programs which differ in array access *patterns* will now be flagged as different. This is a sound approach for difference, as this may include access to illegal bounds (an example for this can be seen in Fig. 16). We note that this definition of description cannot be achieved by instrumenting output or array access through temporary variables and deferring checking their equivalence to the end, as the amount of temporaries needed may be unbound. We therefore require procedures to agree on these locations, in case we wish to verify using the extended definition of difference.

? maybe show a bigger example that better shows how we benefit from partitioning

? talk about how merging early in paths may affect equivalence and differencing behaviors further along

3. Preliminaries

We use the following standard concrete semantics definitions for a program:

Program Location / Label A program location $loc \in Loc$, also referred to as label denoted lab , is a unique identifier for a certain location in a program corresponding to the value of the program counter at a certain point in the execution of the program. We also define two special labels for the start and exit locations of the program as *begin* and *fin* respectively.

Concrete State Given a set of variables Var , a set of possible values for these variables Val and the set of locations Loc , a *concrete program state* is a tuple $\sigma \triangleq \langle loc, values \rangle \in \Sigma$ mapping the set of program variables to their concrete value at a certain program location loc i.e. $values : Var \rightarrow Val$. The set of all possible states of a program P is denoted Σ_P .

Program We describe an imperative program P , as a tuple $(Val, Var, \rightarrow, \Sigma_0)$ where $\rightarrow : \Sigma_P \times \Sigma_P$ is a transition system which given a concrete program state returns the following state in the program and Σ_0 is a set of initial states of the program. Our formal semantics need to deal with errors states therefore we ignore crash states of the programs, as well as inter-procedural programs since our work deals with function calls by inlining and we exclude recursion for now.

Concrete Trace A program trace $\pi \in \Sigma_P^*$, is a sequence of states $\langle \sigma_0, \sigma_1, \dots \rangle$ describing a single execution of the program. Each of the states corresponds to a certain location in the program where the trace originated from. Every program can be described by the set of all possible traces for its run $\llbracket P \rrbracket \subseteq \Sigma_P^*$. We refer to these semantics as concrete state semantics. We also define the following standard operations on traces:

- $label : \Sigma_P \rightarrow Lab$ maps a state to the program label at which it appears.
- $last : \Sigma_P^* \rightarrow \Sigma_P$ returns the last state in a trace.
- $pre : \llbracket P \rrbracket \rightarrow 2^{\Sigma_P^*}$ for a trace π is the set of all prefixes of π .

4. Concrete Semantics

In this section, we define the notion of concrete difference between two programs based on a standard concrete semantics.

4.1 Concrete State Differencing

Comparing two programs P and P' under concrete semantics means comparing their *traces*. Towards that end, we first define the difference between two concrete state.

Intuitively, given two concrete states, the difference between them is the set of variables (and their values) where the two states

map corresponding variables to different values. As variable names may differ between programs, we parameterize the definition by a mapping that establishes a correspondence between variables in P and those in P' . Concrete state differencing only compares values of corresponding variables.

Variable Correspondence A variable correspondence $VC \subseteq Var \times Var'$, is a partial mapping between two sets of program variables, those of P and those of P' . The variable correspondence mapping can be taken as input from the user. However, our experience indicates that is often sufficient to use a standard name-based mapping: $VC_{EQ} \triangleq \{(v, v') | v \in Var \wedge v' \in Var' \wedge name(v) = name(v')\}$.

Concrete State Delta Given two concrete states $\sigma^h \in \Sigma_P, \sigma'^h \in \Sigma_{P'}$, and a variable correspondence mapping VC , we define the concrete state delta $\Delta_S(\sigma^h, \sigma'^h) : \Sigma_P \times \Sigma_{P'} \rightarrow 2^{Var \times Val}$ as the part of the state σ^h where corresponding variables do not agree on values (with respect to σ'^h).

DEFINITION 1. Given two concrete states $\sigma^h \in \Sigma_P, \sigma'^h \in \Sigma_{P'}$, and a correspondence mapping VC , the concrete state delta is defined as:

$$\Delta_S(\sigma^h, \sigma'^h) \triangleq \{(v, val) | (v, v') \in VC \wedge \sigma^h(v) = val \neq \sigma'^h(v')\}$$

Note that Δ_S is not symmetric. In fact, the direction in which Δ_S is used has meaning in the context of a program P and a patched version of it P' . We define $\Delta_S^- = \Delta_S(\sigma^h, \sigma'^h)$ which means the values of the state that was "removed" in P' and $\Delta_S^+ = \Delta_S(\sigma'^h, \sigma^h)$ which stands for the values added in P' . When there is no observable difference between the states we get that $\Delta_S^+(\sigma^h, \sigma'^h) = \Delta_S^-(\sigma^h, \sigma'^h) = \emptyset$, and say that the states are *equivalent*, and write $\sigma \approx_{VC} \sigma'$. When VC is clear from context, we omit the subscript and simply write $\sigma \approx \sigma'$.

EXAMPLE 1. Consider two concrete states $\sigma^h = (x \mapsto 1, y \mapsto 2, z \mapsto 3)$ and $\sigma'^h = (x' \mapsto 0, y' \mapsto 2, w' \mapsto 4)$ and using VC_{EQ} then $\Delta_S^- = (x \mapsto 1)$ since x and x' match and do not agree on value, y and y' agree (thus are not in delta) and z' is not in VC_{EQ} . Similarly, $\Delta_S^+ = (x' \mapsto 0)$.

We now use our notion of concrete state difference to define the difference between concrete program traces. Our goal is to compare traces that are input-equivalent and check whether they are output-equivalent. Therefore, we are interested in comparing traces that originate from equivalent input states. To differentiate traces, we need to compare the states along each trace, but which states should we compare? this is not a trivial question since traces can vary in length and order of states. We need a mapping for choosing the states to be differentiated within the two traces. We define it as following:

Trace Diff Points Given two traces $\pi \in \llbracket P \rrbracket$ and $\pi' \in \llbracket P' \rrbracket$ that originate from equivalent input states, we define a trace index correspondence relation named *trace diff points* denoted DP_π as a matching of indexes specifying states where concrete state delta should be computed. Formally, $DP_\pi \subseteq \{(i, i') | 0 \leq i \leq |\pi|, 0 \leq i' \leq |\pi'|\}$. The question of supplying this matching, in a way that results in meaningful delta, is not a trivial one, we delay this discussion until we define the trace delta.

Trace Delta Now that we have a way of matching states to be compared between two traces, we define the notion of trace differentiation:

DEFINITION 2. Given two traces $\pi \in \llbracket P \rrbracket$ and $\pi' \in \llbracket P' \rrbracket$ that originate from equivalent input states, and a trace index correspondence DP_π we define the trace delta $\Delta_T(\pi, \pi') : \llbracket P \rrbracket \times \llbracket P' \rrbracket \rightarrow$

```
void foo(unsigned x) {      void foo'(unsigned x) {
    unsigned i = 0;          unsigned i = 0;
    lab:guard g = (i >= x);  lab:guard g = (i >= 2*x);
    if (g) return;          if (g) return;
    ...                     ...
    i++;                    i++;
    goto lab;               goto lab;
}
```

Figure 8: P, P' differentiation candidates

$(Loc \rightarrow 2^{Var \times Val})$ as state differentiations between all corresponding states in π and π' .

$$\Delta_T(\pi, \pi') = \{(i, \Delta_S(\sigma_i^h, \sigma'_{i'}^h)) | (i, i') \in DP_\pi, \Delta_S(\sigma_i^h, \sigma'_{i'}^h) \neq \emptyset\}$$

That is, for every $(i, i') \in DP_\pi$ such that $\Delta_S(\sigma_i^h, \sigma'_{i'}^h) \neq \emptyset$, Δ_T will contain the mapping $i \mapsto \Delta_S(\sigma_i^h, \sigma'_{i'}^h)$, thus the result will map certain states in π to their state delta with the corresponding state in π' (deemed interesting by DP_π). Since $\Delta_T(\pi, \pi')$ is based on state difference, we define Δ_T^+ and Δ_T^- similarly to their underlying states difference operations.

One possible choice for DP_π is the endpoints of the two traces $\{(fin, fin')\}$ (assuming they are finite) meaning differentiating the final states of the executions or formally: $\Delta_{fin}^- \triangleq \Delta_{fin}(\pi, \pi') = \{fin \mapsto \Delta_S(\sigma_{fin}^h, \sigma'_{fin'}^h)\}$ (we will also be interested in Δ_{fin}^+). The final state delta may not always be sufficient for truly describing the difference between traces as according to our definition of difference, we would be interested in checking difference at intermediate locations in the program (that emit output or check assertions). Δ_{fin} will only compare final state values and could miss what happened during the execution. This can be overcome by instrumenting the semantics such that the state contains all "temporary" values for a variable along with the trace index (program location is not sufficient here as a trace can loop over a certain location) where the values existed by, for instance, adding temporary variables, allowing for a complete differentiation at the end point. This substantially complicates the selection of VC as it requires relating all of these temporary, indexed, variables. Such a correspondence may be extremely hard to produce. Also the number of variables here can range up to the length of the trace.

Defining the DP over any two traces is a daunting task since traces of separate (although similar) programs can vastly differ.

EXAMPLE 2. Consider the two program versions shown in Fig. ?? and the following traces generated from the input $x = 2$: $\pi = \langle (x \mapsto 2, i \mapsto 0, g \mapsto 0), (x \mapsto 2, i \mapsto 1, g \mapsto 0), (x \mapsto 2, i \mapsto 2, g \mapsto 1) \rangle$ and $\pi' = \langle (x \mapsto 2, i \mapsto 0, g \mapsto 0), (x \mapsto 2, i \mapsto 1, g \mapsto 0), (x \mapsto 2, i \mapsto 2, g \mapsto 0), (x \mapsto 2, i \mapsto 3, g \mapsto 0), (x \mapsto 2, i \mapsto 4, g \mapsto 1) \rangle$, we see that even in this simple program, finding a correlation based on traces alone is hard. However, using program location as a means of correlation, one can produce a meaningful result that describes how the values of i range differently in the new version P' . For example, if we look at all the possible values for i in label lab and differentiate them (as a set) from the values in the patched version (in the same location), we get a meaningful result that i in the patched version can range from $x + 1$ up to $2x$.

4.2 Differencing at Program Labels

We now show that choosing DP based on program locations is quite natural.

Trace Delta using Program Labels Given two traces π, π' and two program labels l, l' we define a trace delta based on all states that are labeled l, l' . First we define π_l as a sub-sequence of π

where only states that are labeled l were chosen ($\pi'_{l'}$ is defined similarly). Next, we denote $\Delta_L(\pi_l, \pi'_{l'})$ as a means for comparing these sequences. As $\pi_l, \pi'_{l'}$ may vary in length and order, we cannot simply define it as applying Δ_S on each pair of states in $(\pi_l, \pi'_{l'})$ by order. In fact, Δ_L can be defined in different way to reflect different concepts of difference, for instance, it can be defined as the differentiating the last states of π_l and $\pi'_{l'}$ (assuming they are both finite) to reflect we are only interested in the final values in that location. We chose to define Δ_L as the difference between the *set of states* which appear in π_l against the set of those in $\pi'_{l'}$. Formally: $\Delta_{L_{set}}(\pi_l, \pi'_{l'}) \triangleq \{\sigma^h \in \text{ran}(\pi_l) \mid \neg \exists \sigma'^h \in \text{ran}(\pi'_{l'}) . \text{s.t.} \Delta_S(\sigma^h, \sigma'^h) = \emptyset\}$. For example consider Fig. ??, for π, π' that originate from $x = 2$ then $\Delta_{L_{set}}(\pi_{lab}, \pi'_{lab'}) = \emptyset$ and $\Delta_{L_{set}}(\pi'_{lab'}, \pi_{lab}) = \{(i' \mapsto 3), (i' \mapsto 4)\}$. We see that this notion of Δ indeed captures a useful description of difference.

The problem of choosing DP is now reduced to the matching of labels as the trace indexing correspondence DP_π defined in Definition 4.1 is induced by the definition over labels. as we need to differentiate sets of states belonging to a certain program label. We require a correspondence of *labels* and therefore we define the label diff points correspondence.

Label Diff Points Given two programs P, P' and their sets of program labels Lab, Lab' , we define a label correspondence relation named label diff points denoted DP_{Lab} as a matching of labels between programs. Formally: $DP_{Lab} \subseteq \{(l, l') \mid l \in Lab, l' \in Lab'\}$. From this point on any mention of the diff-points correspondence DP will refer to label diff-points DP_{Lab} . We address the question of selection of DP_{Lab} in ??.

Now, we will move past the concrete semantics towards *abstract semantics*. This is required as it is unfeasible to describe difference based on traces. However, we need to adjust our concrete semantics before we can correctly define this as the concrete semantics based on individual traces *will not allow us to correlate traces that originate from the same input*. This is the first formal indication of how a separate abstraction, that considers each of the programs by itself, cannot succeed.

4.3 Concrete Correlating Semantics

In this section we shortly define, based on previous definitions, the correlating state and trace which bind the executions of both programs, P and P' , together and we define the notion of delta there. Essentially, these will be states and traces of the product program $P \times P'$ but *only traces that originate from the same input are considered*. This allows us to then define the correlating abstract semantics which is key for successful differencing.

Correlating Concrete State A correlating concrete state $\sigma^\sharp_\times : Var \cup Var' \rightarrow Val$ is a join of concrete state, mapping variables from a pair of programs P and P' to their values. The set of all possible correlating states is denoted $\Sigma^\sharp_{P \times P'}$.

Correlating Concrete Trace A correlating trace π_\times , is a sequence of correlating states $\dots, < \sigma^\sharp_{i_\times}, >, \dots$ describing a dual execution of the two programs where at every point each of the program can perform a single step. The $label_\times$, $last_\times$ and pre_\times operations are defined similarly. We denote by $\llbracket P \times P' \rrbracket$ the set of all traces in $P \times P'$. Again, we restrict to traces that originate from equivalent input state i.e., $\sigma^\sharp_0 \approx_{VC} \sigma'^\sharp_0$.

We must remember however, that the number of traces to be compared is potentially unbounded which means that the delta we compute may be unbounded too. Therefore we must use an abstraction over the concrete semantics that will allow us to represent executions in a bounded way.

5. Abstract Correlating Semantics

In this section, we introduce our correlating abstract domain which allows bounded representation of product program state while maintaining equivalence between correlated variables. This comes at the cost of an acceptable lose of precision of other state information. We represent variable information using standard relational abstract domain. As our analysis is path sensitive, we allow for a set of abstract sub-states, each adhering to a certain path in the product program. This abstraction is similar to the trace partitioning domain as described in [15]. To assure we only consider correlated paths from the product programs (that agree on input), our abstraction will initially assume equality on all inputs. This power-set domain records precise state information but does not scale due to exponential blow-up of number of paths. To reduce state size, we define a special join operation that dynamically partitions the abstract state according to the set of equivalences maintained in each sub-state and joins all sub-states in the same partition together (using the sub-domain lossy join operation). This equivalence criteria allows separation of equivalence preserving paths thus achieving better precision. We start off by abstracting the correlating trace semantics in Sec. 4.3.

In the following, we assume an abstract relational domain $(D^\sharp, \sqsubseteq_D)$ equipped with operations \sqcap_D, \sqcup_D and ∇_D , for representing sets of concrete states in $\Sigma_{P \times P'}$. We separate the set of program variables into original program variables denoted Var (which also include a special added variable for return value, if such exists) and the added guard variables denoted $Guard$ that are used for storing conditional values alone ($Guard$ also include a special added variable for return flag). We assume the abstract values in D^\sharp are constraints over the variables and guards (we denote D^\sharp_{Guard} for abstraction of guards and D^\sharp_{Var} for abstracting original variables), and do not go into further details about the particular abstract domain as it is a parameter of the analysis. We also assume that the sub-domain D^\sharp allows for a sound over-approximation of the concrete semantics (given a sound interpretation of program operations). In our experiments, we use the polyhedra abstract domain [5] and the octagon abstract domain [12].

Correlating Abstract State A correlating abstract program state $\sigma^\sharp \in Lab_\times \rightarrow 2^{D^\sharp_{Guard} \times D^\sharp_{Var}}$, is a set of pairs $\langle ctx, data \rangle \in \Sigma^\sharp$ mapped to a product program label l_\times , where $ctx \in D^\sharp_{Guard}$ is the execution context i.e. an abstraction of guards values via the relational numerical domain and $data \in D^\sharp_{Var}$ also an abstraction of the variables (anything that is not a guard) also using the domain D^\sharp . We separate abstractions over guard variables added by the transformation to GCL and original program variables as there need not be any relationships between guard and regular variables.

Correlating Abstract Semantics Tab. 1 describe the abstract transformers. The table shows the effect of each statement on a given abstract state $\sigma^\sharp = l_\times \mapsto S$. The abstract transformers are defined using the abstract transformers of the underlying abstract domain D^\sharp .

We assume that any program P can be transformed such that it contains only the aforementioned operations (specifically GCL format). We also assume that for $\llbracket g := e \rrbracket^\sharp$ operations, e is a logical operation with binary value. Next, we define the abstraction function $\alpha : 2^{\Sigma^*_{times}} \rightarrow 2^{D^\sharp \times D^\sharp}$ for a set of concrete traces $T \subseteq \Sigma^*_{times}$. As in our domain traces are abstracted together if they share the exact same path, we first define an operation $path : \Sigma^*_{times} \rightarrow Lab^*$ which returns a sequence of labels for a trace's states i.e. what is the path taken by that trace. We also allow applying $path$ on a set of traces to denote the set of paths resulting by applying the function of each of the traces. Finally: $\alpha(T) \triangleq \{\sqcup_{path(\pi)=p} \beta(last(\pi)) \mid p \in path(T)\}$ where

$\llbracket v := e \rrbracket^\#$	$l_x \mapsto \{ \langle ctx, \llbracket v := e \rrbracket_{D^\#}^\#(data) \rangle \mid (ctx, data) \in S \}$
$\llbracket g := e \rrbracket^\#$	$l_x \mapsto \{ \langle \llbracket g := true \rrbracket_{D^\#}^\#(ctx), \llbracket e \rrbracket_{D^\#}^\#(data) \rangle \mid (ctx, data) \in S \} \cup \{ \langle \llbracket g := false \rrbracket_{D^\#}^\#(ctx), \llbracket -e \rrbracket_{D^\#}^\#(data) \rangle \mid (ctx, data) \in S \}$
$\llbracket \text{if}(g)\{s_0\}\text{else}\{s_1\} \rrbracket^\#$	$l_x \mapsto \{ \langle \llbracket g = true \rrbracket_{D^\#}^\#(ctx), \llbracket s_0 \rrbracket_{D^\#}^\#(data) \rangle \mid (ctx, data) \in S \} \cup \{ \langle \llbracket g = false \rrbracket_{D^\#}^\#(ctx), \llbracket s_1 \rrbracket_{D^\#}^\#(data) \rangle \mid (ctx, data) \in S \}$
$\llbracket \text{gotolab} \rrbracket^\#$	$\sigma^\#$

Table 1: Abstract transformers using abstract transformers of the underlying domain $D^\#$. The table describe the effect of each statement on an abstract state $\sigma^\# = l_x \mapsto S$.

```

int f(int x) {          int f'(int x) {
    return x;           return 2*x;
}                      }

```

Figure 9: Single path differentiation candidates

$\beta(\sigma^\#) = \langle \beta_{D^\#}(\sigma^\#|_{Guard}), \beta_{D^\#}(\sigma^\#|_{Var}) \rangle$ i.e. applying the the abstraction function of the abstract sub-domain $\beta_{D^\#}$ on parts of the concrete state applying to *Guards* and *Vars* separately. Our abstraction partitions trace prefixes π by path and abstracts together the concrete states reached by the prefix - *last*(π), using the sub-domain.

Every path in the product program will be represented by a single sub-state of the sub-domain. As a result, all **traces prefixes** that follow the same path to l_x will be abstracted into a single sub-state of the underlying domain. This abstraction fits semantics differencing well, as inputs that follow the same path display the same behavior and will usually either keep or break equivalence together, allowing us to separate them from other behaviors (it is possible for a path to display both behaviors as in Fig. 9 and we will discuss how we are able to manipulate the abstract state and separate equivalent behaviors from ones that offend equivalence). Another issue to be addressed is the fact that our state is still potentially unbound as there may be an infinite number of paths in the program (due to loops).

Partitioning Performing analysis with the powerset domain does not scale as the number of paths in the correlated program may be exponential. We must allow for reduction of state $\sigma^\# = l_x \mapsto S$ with acceptable loss of precision. This reduction via partitioning can be achieved by joining the abstract sub-states in S (using the standard precision losing join of the sub-domain) but to perform this we must first answer the following: (i) which of the sub-states shall be joined together and (ii) at which program locations should the partitioning occur. A trivial partitioning strategy is simply reverting back to the sub-domain by applying the join on all sub-states which may result in unacceptable precision loss as exemplified in Fig. 4. However, by taking a closer look at the final state of the same example:

$$\sigma^\#(fin) = [\langle (g1 = 1, g2' = 0, \equiv_{g1}), (sgn = 1, \equiv_{sgn}) \rangle, \langle (g1 = 0, g2' = 0, \equiv_{g1}), (x < 0, sgn = -1, \equiv_{x,sgn}) \rangle, \langle (g1 = 0, g2' = 1, \equiv_{g1}), (x = 0, sgn = 0, sgn' = 1, \equiv_x) \rangle]$$

One may observe that were we to join the two sub-states that maintain equivalence on $\{x, sgn, g1\}$, it would result in an acceptable loss of precision (of losing the x related constraints). This is achieved by partitioning sub-states according to the set of variables which they preserve equivalence for. This bounds the state size at $2^{|V|}$, where V is the set of correlating variables we wish to track. As mentioned, another key factor in preserving equivalence and maintaining precision is the program location at which the partitioning occurs. The first possibility, which is somewhat symmetric to the first proposed partitioning strategy, is to partition at every join point i.e. after every branch converges. A quick look at Fig. 4's

state after processing the first guarded instruction `if (G1) sgn = -1;`, i.e (we ignored $g2'$ effect at this point for brevity):

$$\sigma^\# = [\langle (g1 = 0, \equiv_{g1}), (x \geq 0, \equiv_{x,sgn}) \rangle, \langle (g1 = 1, \equiv_{g1}), (x < 0, sgn' = -1, \equiv_x) \rangle]$$

suggests that partitioning at join points will perform badly in many scenarios, specifically here as we will lose all data regarding *sgn*. However if we could delay the partitioning to a point where the two programs "converge" (after the following `if (G1') sgn' = -1;` line), we will get a more precise temporary result which preserves equivalence. Therefore we define another partitioning points, which are defined during the correlating program construction process described in Section 6, and are basically locations where the two (guarded) programs have syntactically converged (we found two lines that are syntactically equal, sans tagging). Another possible location for differencing are our *differencing points*. Partitioning at these locations is essentially more precise than at correlation points. We remind that diff-points are product program locations where both programs conceptually converge as they are about to emit output i.e. both programs have finished "preparing" the next portion of output therefore, if equivalence exists - it must exist now, therefore this is a prime candidate for a partitioning point. **Our evaluation includes applying each of our strategies along with each of the points. Intuitively, the results should range from least precise using the <Join-All,At-Join> strategy and point to most precise in the <Join-Equiv,At-Diff> scenario and this is indeed the case as we show in Section 7 (not taking into account the no-partitioning scenario which is naturally most precise).**

Widening In order for our analysis to handle loops we require a means for reaching a fixed point. As our analysis advances over a loop and state is transformed, it may keep changing and never converge unless we apply the widening operator to further over-approximate the looping state and arrive at a fixed point. We have the widening operator of our sub-domain at our disposable, but again we are faced with the question of how we apply this operator, i.e. which pairs of sub-states $\langle ctx, data \rangle$ from $\sigma^\#$ should be widened with which. A first viable strategy, similar to the first partitioning strategy, is to perform an overall join operation on all pairs which will result in a single pair of sub-states and then simply apply the widening to this sub-state using the sub-domain's ∇ operator. If we examine applying this strategy to $sum \bowtie sum$ from Fig. 6, we see that it will successfully arrive at a fixed point that also maintains equivalence as all sub-states maintain equivalence at loop back-edges. Now let us try applying the strategy to the more complex $sum \bowtie sum'$ as seen in Fig. 6. First we mention that as sum' introduces a return statement under the $len > max$ condition, the example shows an extra r' guard for representing a return (this exists in all GCL programs but we omitted it so far for brevity). While analyzing, once we pass that first conditional, our state is split to reflect the return effect:

$$\sigma^\# = [c_1 = \langle (r' = 1), (len = len' > max, return_value' = -1, result = result' = 0) \rangle, c_2 = \langle (r' = 0), (len = len' < max, result = result' = 0) \rangle]$$

As we further advance into the loop, c_2 will maintain equivalence

```

unsigned max = ...;
int sum'(int arr[], unsigned len) {
    int result = 0;
    if (len > max)
        return -1;
    for (unsigned i = 0; i < len; i++)
        result += arr[i];
    return result;
}

unsigned max' = ...;
int sum(int arr[], unsigned len) {
    unsigned len' = len;
    int arr'[] = arr;
    int result = 0;
    int result' = 0;
    guard r' = (len' > max');
    if (r') return_value' = -1;
    if (r') r' = 0;
    {
        unsigned i = 0;
        unsigned i' = 0;
    1: guard g = (i < len);
    1': guard g' = (i' < len');
        if (g) result += arr[i];
        if (r') if (g') result' += arr'[i'];
        if (g) i++;
        if (r') if (g') i'++;
        if (g) goto 1;
        if (r') if (g') goto 1';
    }
}

```

Figure 10: Patched sum' and correlating $sum \bowtie sum'$

but c_1 will continue to update the part of the state regarding untagged variables (since $r' = 1$ in c_1 and it will not consider any of the commands guarded by r'), specifically it will change $result$ continuously, preventing the analysis from reaching fixed point. We would require widening here but using the naive strategy of a complete join will result in aggressive loss of precision, specifically losing all information regarding $result$. The problem originates from the fact that prior to widening, we joined sub-states which adhere to two different loop behaviors: one where both sum and sum' loop together (that originated from $len < max$) and the other where sum' has existed but sum continues to loop ($len \geq max$). Ideally, we would like to match these two behaviors and widen them accordingly. We devised a widening strategy that allows to do this as it basically matches sub-states that adhere to the same behavior, or loop-paths. We do this by using *guards* for the matching. If two sub-states agree on their set of guards, it means they represent the same loop path and can be widened as the latter originated from the former (widening operates on subsequent iterations). In our example, using this strategy will allow the correct matching of states after consequent $k, k + 1$ loop iterations:

$$\begin{aligned}
& \sigma_k^\# = [\\
& \quad c_1 = \langle (r' = 1, g=1, g'=0), \\
& \quad (len = len' > max, return_value' = -1, result' = 0, result = \sum_{i=0}^k arr[i]) \rangle \\
& \quad c_2 = \langle (r' = 0, g=1, g'=1), \\
& \quad (len = len' < max, result = result' = \sum_{i=0}^k arr[i]) \rangle \\
&]
\end{aligned}$$

And:

$$\begin{aligned}
& \sigma_{k+1}^\# = [\\
& \quad c_1 = \langle (r' = 1, g=1, g'=0), \\
& \quad (len = len' > max, return_value' = -1, result' = 0, result = \sum_{i=0}^{k+1} arr[i]) \rangle \\
& \quad c_2 = \langle (r' = 0, g=1, g'=1), \\
& \quad (len = len' < max, result = result' = \sum_{i=0}^{k+1} arr[i]) \rangle \\
&]
\end{aligned}$$

As we can identify the states predecessors by simply matching the guards. c_1 will be widened for a precise description of the

difference shown as $\langle len = len' > max, return_value' = -1, return_value' = \top \rangle$.

5.1 Correlating Abstract State Differencing

Given a state in our correlating domain, we want to determine whether equivalence is kept and if so under which conditions it is kept (for partial equivalence) or determine there is difference and characterize it. As our state may hold several pairs of sub-states, each holding different equivalence data, we can provide a verbose answer regarding whether equivalence holds. We partition our sub-states according to the set of variables they hold equivalence for and report the state for each equivalence partition class. Since we instrument our correlating program to preserve initial input values, for some of these states we will also be able to report input constraints thus informing the user of the input ranges that maintain equivalence. In the cases where equivalence could not be proved, we report the offending states and apply a differencing algorithm for extraction of the delta. Fig. 9 shows an example of where our analysis is unable to prove equivalence (as it is sound), although part of the state does maintain equivalence (specifically for $x = 0$). This is due to the abstraction being too coarse. We describe an algorithm that given a sub-state $d \in D^\#$, computes the differentiating part of the sub-state (where correlated variables disagree on values) by splitting it into parts according to equivalence. This is done by treating the relational constraints in our domain as geometrical objects and formulating delta based on that.

Correlating Abstract State Delta Given a sub-state d and a correspondence VC , the correlating state delta $\Delta_A(d)$, computes abstract state differentiation over d . The result is an abstract state $\sqsubseteq rd$ approximating all concrete values for variables correlated by VC , that differ between P and P' . Formally, the delta is simply the abstraction of the concrete trace deltas $\alpha(\cup_{path} \Delta_T^+)$, $\alpha(\cup_{path} \Delta_T^-)$ where deltas are grouped together by path and then abstracted. But it is not clear as to how we compute this set differencing on the correlating abstraction. Instead, we consider the geometric representation of the domain and applied operations there for the extraction of delta, as following:

1. d_{\equiv} is a state abstracting the concrete states shared by the original and patched program. It is achieved by computing: $d_{\equiv} \triangleq d|_{V=V'} \equiv d \sqcap \bigwedge \{v = v' | VC(v) = v'\}$.
2. $\overline{d_{\equiv}}$ is the negated state i.e. $D^\# \setminus d_{\equiv}$ and it is computed by negating d_{\equiv} (as mentioned before, all logical operations, including negation, are defined on our representation of an abstract state).
3. Eventually: $\Delta_A(d) \triangleq d \sqcap \overline{d_{\equiv}}$ abstracts all states in $P \times P'$ that where correlated variables values do not match.
4. $\Delta_A(d)^+ = \Delta_A(d)|_{V'}$ is a projection of the differentiation to display values of P' alone i.e. "added values".
5. $\Delta_A(d)^- = \Delta_A(d)|_V$ is a projection of the differentiation to display values of P alone i.e. "removed values".

Applying the algorithm on Fig. 9's P and P' where $rd = \{retVal' = 2retVal\}$ will result in the following:

1. $d_{\equiv} = \{retVal' = 0, retVal = 0\}$.
2. $\overline{d_{\equiv}} = [\{retVal' > 0\}, \{retVal' < 0\}, \{retVal > 0\}, \{retVal < 0\}]$
3. $\Delta_A(d) = [\{retVal' = 2retVal, retVal' > 0\}, \{retVal' = 2retVal, retVal' < 0\}, \{retVal' = 2retVal, retVal > 0\}, \{retVal' = 2retVal, retVal < 0\}]$
4. $\Delta_A(d)^+ = [\{retVal' > 0\}, \{retVal' < 0\}]$
5. $\Delta_A(d)^- = [\{retVal > 0\}, \{retVal < 0\}]$

Figure 11: Delta computation geometrical representation.

We see that displaying the result in the form of projections is ill-advised as in some states differentiation data is represented by relationships on correlated variables alone, thus projecting will lose all data and we will be left with a less informative result. A geometrical representation of Δ_A calculation can be seen in Fig. 11.

From this point forward any mention of 'delta' (denoted Δ) will refer to the correlating abstract state delta (denoted Δ_A). We claim that Δ is a correct abstraction for the concrete state delta which allows for a scalable representation of difference we aim to capture.

6. Correlating Program

In this section we will describe how we construct our correlating program $P \bowtie P'$. The process of correlating attempts to find a better interleaving of programs for a more precise differentiation. The building process also instruments $P \bowtie P'$ with the required differencing points DP which allow reporting of difference and finds correlation points CP which define the locations for out partitioning. We also allow a user defined selection of DP and CP . We start off with a discussion of differencing points and how the analysis calculates difference

6.1 Construction of $P \bowtie P'$

The idea of a correlating program is similar to that of self-composition [1, 17], but the way in which statements in the correlating program are combined is carefully designed to keep the steps of the two programs close to each other. Rather than having the patched program sequentially composed after the original program, our correlating program interleaves the two versions. Analysis of the correlating program can then recover equivalence between values of correlated variables even when equivalence is *temporarily* violated by an update in one version, as the corresponding update in the other version follows shortly thereafter.

We will generally describe the process of constructing the correlating program. The correlating program is an optimized reduction over $P \times P'$ where not all pairs of (σ, σ') are considered, but only pairs that result from a controlled execution, where correlating instructions in P and P' will execute adjacently. This will allow for superior precision.

The input for the correlation process are two program (P, P') in C. The first step involves transforming both programs to a normalized guarded instruction form (P^G, P'^G) . Next, a vector of *imperative commands* I (and I' respectively) is extracted from each program for the purposes of performing the syntactic diff. An imperative command in our GCL format is defined to be either one of $v := e | gotol f(\dots)$ as they effectively change the program state (variable values, excluding guard). Then a syntactical diff [9] is computed over the vectors (I, I') . One of the inputs to the diff process is VC as it is needed to identify correlated variables and the diff comparison will regard commands differing by variable names which are correlated by VC as equal. The result of the last step will be a vector I_{diff} specifying for each command in I, I' whether it an added command in P' (for I') marked $+$, a deleted command from P (for I) marked $-$, or a command existing in both versions $=$. This diff determines the order in which the commands will be interleaved in the resulting $P \bowtie P'$ as we will iterate over the result vector I_{diff} and use it to construct the correlating program. We remind that since I, I' are only the imperative commands, we cannot use it directly as $P \bowtie P'$. Instead we will use the imperative

commands as markers, specifying which chunk of program from P or P' should be taken next and put in the result. The construction goes as follows: iterate over I_{diff} and for every command c (c') labeled l_c ($l_{c'}$):

- read P (P') up to label l_c ($l_{c'}$) including into block B_c ($B_{c'}$)
- for $B_{c'}$, tag all variables in the block.
- emit the block to the output.
- delete B_c ($B_{c'}$) from P (P').

The construction is now complete. We only add that at the start of the process, we strip P' of its prototype and add declarations for the tagged input variables, initializing them to the untagged version. As mentioned CP is also a product of the construction, and it's defined using $=$ commands: after two $=$ commands are emitted to the output, we add an instrumentation line, telling the analysis of the correlation point. One final observation regarding the correlating program is that it is a legitimate program that can be run to achieve the effect of running both versions. We plan to leverage this ability to use dynamic analysis and testing techniques such as fuzzing ?? and directed automated testing ?? on the correlating program in our future work.

7. Evaluation

We mainly tested our tool on the GNU core utilities, differencing versions 6.10 and 6.11. Our benchmark is composed of ? patched programs which include changes to numerical variables. Our tool was able to precisely describe the difference. We also tested our tool on a few handpicked patches taken from the Linux kernel and the Mozilla Firefox web browser.

We implemented a correlating compiler named CCC which creates correlating programs from any two C programs as well as a differencing oriented dataflow analysis solver for analyzing correlated programs, both tools use the LLVM and CLang compiler infrastructure. We analyze C code directly since it is more structured, has type information and keeps a low number of variables, as opposed to intermediate representation. We also benefit from our delta being computed over original variables. As mentioned in Section 6, we normalize the input programs before unifying them for a simpler analysis. We employed the APRON abstract numerical domain library and conducted our experiments using several domains including interval, octagon and polyhedra. We found that octagon and polyhedra domain provide the same level of precision in our experiments (while interval produces many false positives). **We do inlining manually. We assume syscalls do not affect locals**

All of our experiments were conducted running on a Intel(R) Core-i7(TM) processor with 4GB. The results we present have been trimmed down for brevity, the full version of the results can be found at [link to full results](#)

7.1 Results

Producing delta from abstract state Fig. 12 depicts a patch made to the net/sunrpc/addr.c module in the Linux kernel SUNRPC implementation v2.6.32-rc6 which is aimed at removing an off-by-two array access out of bounds error in the original program. Analyzing the correlated program for the two versions will produce the following result:

σ_1 :	
	returned' = true
	returned = false
	uaddr_len' ≤ RPCBIND_MAXUADDRLEN - 2
	uaddr_len' ≥ RPCBIND_MAXUADDRLEN

The difference is captured by one sub-state where the execution has ended in the patched program but continues in the original. We instrument the correlating program with a return flag to capture the

```

size_t rpc_uaddr2sockaddr (const size_t uaddr_len, ...) {
    char buf[ RPCBIND_MAXUADDRLEN ];
    ...
    if ( uaddr_len > sizeof ( buf ) )
        return 0;
    ...
    (*)_1
    buf [ uaddr_len ] = '\n';
    buf [ uaddr_len + 1 ] = '\0';
    ...
}

```

Original

```

size_t rpc_uaddr2sockaddr (const size_t uaddr_len, ...) {
    char buf[ RPCBIND_MAXUADDRLEN ];
    ...
    if ( uaddr_len > sizeof ( buf ) - 2 )
        return 0;
    ...
    (*)_1
    buf [ uaddr_len ] = '\n';
    buf [ uaddr_len + 1 ] = '\0';
    ...
}

```

Patched

Figure 12: Original and patched version of Linux kernel net/sunrpc/addr.c v2.6.32-rc6 module

```

nsresult SetTextInternal (int textLength, int aCount,
                        int aLength, int aOffset,
                        PRUnichar * aBuffer) {
    PRInt32 newLength = textLength - aCount + aLength ;
    PRUnichar * to;
    ...
    (*)_1
    memcpy (to + aOffset , aBuffer ,
            aLength * sizeof ( PRUnichar ));
    ...
}

```

Original

```

nsresult SetTextInternal (int textLength, int aCount,
                        int aLength, int aOffset,
                        PRUnichar * aBuffer) {
    PRInt32 newLength = textLength - aCount + aLength ;
    PRUnichar * to;
    ...
    if ((unsigned)newLength > (1 << 29))
        return NS_ERROR_DOM_DOMSTRING_SIZE_ERR;
    (*)_1
    memcpy (to + aOffset , aBuffer ,
            aLength * sizeof ( PRUnichar ));
    ...
}

```

Patched

Figure 13: Original and patched version of Mozilla Firefox /content/base/src/nsGenericDOMDataNode.cpp v3.5-3.6.4 module

difference as otherwise equivalence holds (none of the original variables change).

Another example, taken from CVE-2010-1196 advisory regarding Firefox’s heap buffer overflow on 64-bit systems is shown in Fig. 13 (vulnerable part of the function only). Firefox 3.5 and 3.6 (upto 3.6.4) contain a heap buffer overflow vulnerability which is caused by an integer overflow. Due to the amount of data needed to trigger the vulnerability (> 8GB), this is only exploitable on 64-bit systems. The vulnerable code is found in /content/base/src/nsGenericDOMDataNode.cpp of the Mozilla code base and was adapted to C for analysis purposes.

Here, we need to describe a more complex and non-convex constraint that leads to difference. Running DIZY with partitioning would produce the following result for this patch:

σ_1 :
returned' = true
returned = false
return value' = NS_ERROR

The difference in state is described correctly as indeed the only change in values for the patch scenario would be the return value and the early return of the patched version. However, we did not preserve the conditional constraints as they are non-convex. Running the same analysis with no partitioning (this is feasible as the procedure does not loop) will produce:

σ_1 :	σ_2 :
returned' = true	returned' = true
returned = false	returned = false
newLength > 536870912	newLength > -3758096384
	newLength < 0
return value' = NS_ERROR	return value' = NS_ERROR

Now we see that the (unsigned)newLength > (1 << 29) constraint has been successfully encoded in two offending states, each holding a part of the problematic range.

Capturing complex delta Fig. 14 depicts a patch made to the char_to_clump function in version 6.11. The patch replacing the execution of the line `input_position += width;`, which originally executed unconditionally, with a conditional structure that in the new version, allows the line to execute only under certain complex conditions. Since the variables handled in this patch (the global `input_position` and return value `chars`) emit output, describing how their values changed and under which terms is important, especially as the patch cannot be easily parsed by a programmer to understand its meaning. Our analysis produces the following description for the differencing point at the return point:

σ_1 :	σ_2 :	σ_3 :
chars' = 0	input_position' = 0	input_position' = 0
input_position = width	input_position < width	input_position > width
input_position < 0	width < 0	input_position ≤ 0
input_position' = 0		

The result convey the difference in the output variable values alongside some of conditions under which the difference occurs. The result is composed of three sub-states featuring difference (the offending variables appear before each sub-state) and adhere to two added paths in the patched program. The first sub-state belongs to the first branch in the added conditional: the difference is comprised of (i) the new value of `input_position` is 0 as opposed to it being `width` in the former version (the analysis took the `input_position += width;` line into account and incorporated knowing that `input_position = 0` from the branch condition). The analysis also deduced that the old `input_position` is negative under the same input as the branch condition dictates that `width` is negative. (ii) `chars` in the new program is 0 under this path. The two other sub-states adhere to the second added path in the conditional and track a difference for `input_position` alone, basically stating that `input_position` under this path used to assume values in ranges $[-\text{inf}, \text{width}]$ and $[\text{width}, 0]$ but now is simply 0. The splitting of this path into two cases is a result of expressing the non-convex `input_position ≠ 0` condition from the first branch conditional using two sub-states. Running the analysis again while saving the initial values of variables and parameters will produce an even more precise result as following:

```

int input_position;

bool char_to_clump(char c) {
    int width;
    ...
    input_position += width;
    (*)1
    ...
    return chars;
}

```

coreutils pr.c v6.10

```

int input_position;

bool char_to_clump'(char c) {
    int width;
    ...
    if (width < 0 && input_position == 0) {
        chars = 0;
        input_position = 0;
    } else if (width < 0 && input_position <= -width) {
        input_position = 0;
    } else {
        input_position += width;
    }
    (*)1
    ...
    return chars;
}

```

coreutils pr.c v6.11

Figure 14: Original and patched version of coreutils pr.c's char_to_clump procedure

σ_1 :	σ_2 :	σ_3 :	partitioning:	
input_position ₀ = 0 chars' = 0 input_position = width input_position < 0 input_position' = 0	input_position ₀ < -width input_position ₀ < 0 input_position' = 0 input_position < width width < 0	input_position ₀ < -width input_position ₀ > 0 input_position' = 0 input_position > width input_position <= 0	σ_1 : s_len = 0 i' = i s_len' - 1 ≥ i'	σ_2 (equivalent): s_len' = s_len i' = i s_len' - 1 ≥ i'
This result describes constraints on the procedure's input under which the difference exist. Another product of the analysis, which we do not show here, are sub-states describing paths which the patch did not affect.				
<div> $(*)_2$: σ_1 (equivalent): s_len' = s_len i' = i s_len' - 1 ≥ i' </div>				

Maintaining Equivalence and Reporting Difference in Loops

Fig. 15 shows two version of the java logicalValue() method taken from [?], adapted to C. This example features semantic preserving refactoring modification (introducing the elapsed variable and THRESHOLD constant, simplifying a conditional and moving the return statement out of branch block) and one semantic change where 1 is returned instead of old in case $currentTime - t < 100$. The challenge in this example is of course proving equivalence over the loop branch and reporting difference for the negated path. Using a separate analysis, we would have to deduce at the following loop invariant: $val = \sum_{i=0}^{data.length} data[i]$ in order to show equality. However, as our abstraction focuses on variable relationships and our correlating program allows us to interleave the two loops in lock-step, all our analysis needs to deduce is the $val = valw$ constraint. As we apply widening to converge, the constraint will be kept, allowing us to establish equivalence for the looping path. Thus DIZY will report the following differencing state for the exit point of logicalValue():

```

 $\sigma_1$ :
currentTime - t < 100
return value = old
return value' = 1

```

We note that the example was run in [?] by unrolling 2 steps of the loop. Next we shall explore a different loop scenario where all paths in the programs contain loops and only some of them maintain equivalence.

Fig. 16 shows part of coreutils md5sum.c bsd_split_3 function that was patched in version version 6.11 to disallow 0-length inputs. Although this patch seems trivial, analyzing it is challenging as it affects the behavior of loops i.e. unbound path lengths. The main challenge in this example, is separating the path where s_len is 0, which results in the loop index i ranging within negative values (which result in array access out of bounds fault), from the rest of the behaviors that maintain equivalence, throughout the widening process which is required for the analysis to reach a fixed point. The result of our analysis is as follows, per differencing point, with

We can see the analysis successfully reports a difference for the singularity point $s_len = 0$ inside the loop, precisely describing the scenario where i' is negative. We can also see the other equivalent state existing within the loop which depicts the results of the widened analysis for all other paths (the $s_len \neq 0$ constraint is not existing there due to partitioning as we will soon show). The differencing sub-state will be omitted once we move past the loop as the $i \leq -1$ constraint will not allow it to exist beyond the loop body (**we do not account for overflow at this point in our work as ???**) thus we are left with the equivalent state alone after the loop which correctly expresses the fact that the programs are equivalent at this point (since both i's converged at 0). We can see that the result at the second differencing point has lost precision since it does not reflect the $i = 0$ constraint. The loss of this constraint is, again, due to partitioning as both sub-states that describe exiting the loop and the one describing entering the loop, hold equivalence for all variables and are joined to together and lose the extra constraint information. **? describe the canonization for the example** If we analyze the same example with no partitioning we get:

$(*)_1$:	σ_2 (equivalent):	σ_3 (equivalent):
σ_1 : s_len = 0 s_len' = 0 i ≤ -1	s_len' = s_len i' = i s_len' - 1 ≥ i'	s_len' = s_len i = 0 s_len' ≥ 1
$(*)_2$:		
σ_1 (equivalent): s_len' = s_len i' = i i = 0 s_len' ≥ 1		

Which further separates the paths in the program, allowing for a different sub-state for the $i = 0$ and $i \neq 0$ substates (again, the $i \neq 0$ constraints was lost when joining together the $i > 0$ and $i < 0$ states as they both adhere to the same path and hold the same guard values). This extra precision is beneficial, but we still managed to

```

int logicalValue(int t) {
    if (!(currentTime - t >= 100)) {
        return old;
    } else {
        int val = 0;
        for (int i = 0 ; i < data.length; i++)
            val = val + data[i];
        old = val;
        return val;
    }
}

const int THRESHOLD = 100;
int logicalValue(int t) {
    int elapsed = currentTime - t;
    int val = 0;
    if (elapsed < THRESHOLD) {
        val = 1;
    } else {
        for (int i = 0 ; i < data.length; i++)
            val = val + data[i];
        old = val;
    }
    return val;
}

```

Figure 15: Two versions of the `logicalValue()` procedure taken from [?]

```

bool bsd_split_3 (char *s, size_t s_len,...) {
    int i = s_len;
    i--;
    while (i && s[i] != '\0') { (*)_1
        i--;
    }
    ...
    (*)_2
}

bool bsd_split_3 (char *s, size_t s_len,...) {
    int i = s_len;
    i--;
    if (s_len == 0) return false;
    i = s_len - 1;
    while (i && s[i] != '\0') { (*)_1
        i--;
    }
    ...
    (*)_2
}

```

coreutils md5sum.c v6.10 coreutils md5sum.c v6.11

Figure 16: Original and patched version of `coreutils md5sum.c`'s `bsd_split_3` procedure

supply a satisfactory result using the more scalable partitioning by equivalence technique.

8. Related Work

Our work has been mainly inspired by recent work identifying program differencing as having vast security implications [2, 16] as well as advancements made in the field of under-approximations of program equivalence [6, 11, 13, 14].

We rely on classic methods of abstract interpretation [4] for presenting an over approximating solution for semantic differencing and equivalence. To achieve this we devised a static analysis over a newly defined construct we call a correlating program. The idea of a correlating program is similar to that of self-composition [1, 17] except that we compose two different programs in a interleaving designed to maintain a close correlation between them. The use of a correlating construct for differencing is novel as previous methods mainly use sequential composition [6, 11, 13, 14], disregarding possible program correlation.

We base our analysis on a relational abstraction [5, 12] that allows us to reason about variables of different programs. The abstraction is further refined towards a disjunctive domain, similar to trace partitioning [15] and we use an equivalence based partitioning criteria, which is apt to our purposes.

Symbolic execution based methods [13, 14] offer practical equivalence verification techniques for loop and recursion free programs, with bound state space.

[?]

Determining corresponding components As suggested in [?], one possibility is to rely on the editing sequence that creates the new version from the original one. Another option is using various syntactic differencing algorithms as a base for computing correspondence tags.

their idea for computing correspondence, is to minimize the “size of change”. They have two different notions of size of change.

[?] introduced a correlating heap semantics for verifying linearizability of concurrent programs. In their work, a correlating heap semantics is used to establish correspondence between a concurrent program and a sequential version of the program at specific linearization points.

A. Worklist

A.1 Points to Hammer

1. a special kind of self composition, where correlated steps are kept together. This is particularly important when handling loops.

A.2 TODO

1. run on uc-kee examples.
2. Consider describing each sub-state as a single (possibly looping) path of execution in both programs that originated from the same input.

A.3 Questions

1. how come you don't need the “product program”?
2. what are the theorems that you provide? (no reason to have definitions if there are no theorems).
3. what abstract domains can we use as “underlying domains” for our abstraction? Do we have any particular requirements from the abstract domains (one requirement is being relational).
4. what makes a “patched version of a program” different from just saying “a different program”? In other words - what are the requirements on the difference between P and P' ?
5. can we claim that our abstraction “forgets” paths along which equivalence is established, but keeps apart paths along which the is a difference, hoping that it will re-converge later?
6. (intuition only) what if we correlate badly and lose soundness? one can propose a 2 "trick" programs that correlating them our way gives a result that loses difference.

```

if (input % 2 == 0) goto 2 else goto 4;
s := input+2;
goto 5;
s := input+3;
if (s>input) goto 6 else goto ERROR;
ptr := realloc(ptr,s);
// use ptr[0], ptr[1], ... ptr[input-1]

```

Figure 17: Example from [?]]

7. WHY DO WE CHECK DIFFERENCE ABOVE THE SUB-STATE LEVEL? doesn't that mean we compare different paths? isn't that bad?

A.4 Useful text fragments

Differential static analysis is useful for regression debugging [?], and may also lead to an automated approach for patch-based exploit generation [?]. Such automated exploit generation enables the organization releasing a patch to estimate the attack surface exposed by its release. Furthermore, it enables the organization releasing the patch to reduce vulnerability to manual and automated patch-based exploitation.

References

- [1] BARTHE, G., D'ARGENIO, P. R., AND REZK, T. Secure information flow by self-composition. In *Proceedings of the 17th IEEE workshop on Computer Security Foundations* (Washington, DC, USA, 2004), CSFW '04, IEEE Computer Society, pp. 100–.
- [2] BRUMLEY, D., POOSANKAM, P., SONG, D., AND ZHENG, J. Automatic patch-based exploit generation is possible: Techniques and implications. In *Proceedings of the 2008 IEEE Symposium on Security and Privacy* (Washington, DC, USA, 2008), IEEE Computer Society, pp. 143–157.
- [3] CHAKI, S., GURFINKEL, A., AND STRICHMAN, O. Regression verification for multi-threaded programs. In *Verification, Model Checking, and Abstract Interpretation*, V. Kuncak and A. Rybalchenko, Eds., vol. 7148 of *Lecture Notes in Computer Science*. Springer Berlin / Heidelberg, 2012, pp. 119–135.
- [4] COUSOT, P., AND COUSOT, R. Abstract interpretation: A unified lattice model for static analysis of programs by construction of approximation of fixed points. In *POPL* (1977), pp. 238–252.
- [5] COUSOT, P., AND HALBWACHS, N. Automatic discovery of linear restraints among variables of a program. In *Conference Record of the Fifth Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (Tucson, Arizona, 1978), ACM Press, New York, NY, pp. 84–97.
- [6] GODLIN, B., AND STRICHMAN, O. Regression verification. In *DAC* (2009), pp. 466–471.
- [7] HORWITZ, S. Identifying the semantic and textual differences between two versions of a program. In *Proceedings of the ACM SIGPLAN 1990 conference on Programming language design and implementation* (New York, NY, USA, 1990), PLDI '90, ACM, pp. 234–245.
- [8] HORWITZ, S., PRINS, J., AND REPS, T. Integrating noninterfering versions of programs. *ACM Trans. Program. Lang. Syst.* 11, 3 (July 1989), 345–387.
- [9] HUNT, J. W., AND MCILROY, M. D. An algorithm for differential file comparison. Tech. rep., Bell Laboratories, 1975.
- [10] JIN, W., ORSO, A., AND XIE, T. Bert: a tool for behavioral regression testing. In *Proceedings of the eighteenth ACM SIGSOFT international symposium on Foundations of software engineering* (New York, NY, USA, 2010), FSE '10, ACM, pp. 361–362.
- [11] KAWAGUCHI, M., LAHIRI, S. K., AND REBELO, H. Conditional equivalence. Tech. rep., MSR, 2010.
- [12] MINÉ, A. The octagon abstract domain. *CoRR abs/cs/0703084* (2007).
- [13] PERSON, S., DWYER, M. B., ELBAUM, S. G., AND PASAREANU, C. S. Differential symbolic execution. In *SIGSOFT FSE* (2008), pp. 226–237.
- [14] RAMOS, D., AND ENGLER, D. Practical, low-effort equivalence verification of real code. In *Computer Aided Verification*, vol. 6806 of *LNCIS*. Springer, 2011, pp. 669–685.
- [15] RIVAL, X., AND MAUBORGNE, L. The trace partitioning abstract domain. *ACM Trans. Program. Lang. Syst.* 29, 5 (Aug. 2007).
- [16] SONG, Y., ZHANG, Y., AND SUN, Y. Automatic vulnerability locating in binary patches. In *International Conference on Computational Intelligence and Security, CIS 2009* (2009), pp. 474–477.
- [17] TERAUCHI, T., AND AIKEN, A. Secure information flow as a safety problem. In *Proceedings of the 12th international conference on Static Analysis* (Berlin, Heidelberg, 2005), SAS'05, Springer-Verlag, pp. 352–367.