

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 original program and its patched version. Towards that end, we first construct a *union program* in which these relationships can be identified, and then use a *correlating abstract domain* to compute a sound approximation of these relationships. To 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 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](#). We evaluate over 50 patches and show that for these programs, the tool often manages to establish equivalence, reports useful approximation of semantic differences when differences exist, and reports only a few false differences.

1. Introduction

When applying a patch to a procedure, the programmer has very limited means for acquiring a description of the change the patch made to the procedure behavior.

Given a program P and a patched version of the program P' our goal is to determine the difference between $\llbracket P \rrbracket$ and $\llbracket P' \rrbracket$.

Existing Techniques Existing techniques for proving patch equivalence will only supply the programmer with a binary answer [4] as to the (input-output) equivalence a program but no description of the difference is supplied. Further work [6] allows refining the equivalence proof by providing a set of constraints under which equivalence is desired but requires the programmer to manually deduce these. Other techniques for describing difference [?] which rely on symbolic execution supply unsound results as they are limited by loops and essentially cover a subset of program behavior. We present a novel approach which allows for a sound description of difference for programs with loops. Our technique employs methods of abstract interpretation for over-approximating the difference in behaviors, by focusing the abstraction on the *relationships* between program behaviors i.e. between variables values (data) and conditionals (path) in the two versions.

In contrast to existing techniques, our approach allows checking for equivalence in every point of execution, and for every variable, while previous approaches focus only on input-output equivalence.

This enables detection of key differences that impact the correctness of the patch: if the changed behavior includes a bug manifested by a local variable (for instance: array index out of bounds), we will detect and describe it while previous work only detected it when propagated to the output and equivalence may have been reported although a bug was introduced. This also provides a challenge as we need to carefully choose the program locations where we check for difference otherwise we will spuriously detect difference.

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 resulting program. Another advantage of this new construct, is that you may apply other methods for equivalence checking directly on it [9] as the correlation allows for a much more fine-grained equivalence checking (between local variables and not only output).

Correlating Abstractions 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 out analysis over the correlating program while updated the domain to reflect program behavior.

To establish equivalence between correlated variables and precisely capture differences, our domain has to maintain correlating information even when other information is abstracted away.

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. To avoid exponential blowup, our join operator may over-approximate numerical information as long as equivalence between correlated variables is preserved.

In some cases, it would have been sufficient to use alternative domains that are capable of representing richer information, such as interval polyhedra [3], or other numerical domains that can represent non-convex information (e.g., [?]). The recent donut domain [?] may be of particular interest for this purpose. However, the general principle of having to preserve correlating information even when information about the values is abstracted away, holds in all of these cases.

In this paper, we present a technique based on abstract interpretation, that is able to compute an over-approximation of the difference between numerical programs or establish their equivalence when no difference exists. The approach is based on two key ideas: (i) create a *union program* that captures the behavior of both the original program and its patched version; (ii) analyze the union pro-

gram with a *correlating domain* that captures relationships between values of variables in the original program and values of variables in the patched version.

The idea of a union program is similar to that of self-composition [24, 10], but the way in which statements in the union 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 union program interleaves the two versions. Analysis of the union 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.

check whether Aiken paper SAS'05 does some sort of interleaving as well

also need to say that there is the problem of choosing differencing points

1.1 Main Contributions

The main contributions of this paper are as follows:

- we phrase the problem of semantic differential analysis as an analysis of a union program — a single program that represents an original program and its patched version.
- we present an approach for analyzing differences over the union program using a correlating abstract domain. Our approach is sound — if there is a difference at a differentiation point, we cannot miss it. However, since we over-approximate differences, our approach may report false differences due to approximation.
- 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 over 50 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.

mention classical work like [?] [?]

2. Overview

In this section, we informally describe our approach with a simple example program.

2.1 Motivating Example

Fig. 2.1 shows a looping procedure and its patched version, including two lines of patch where patch (1) is a bug fix and patch (2) is refactoring. We aim to find semantic differences between these versions which optimally would be the following observations:

- when entering the loop, variables `s_len`, `i` can no longer be 0, `-1` respectively.
- inside the loop, `i` can no longer range between $-\infty$ and -1 .
- the second line of patch makes no difference.

One might consider performing a separate analysis on each of these procedures and comparing the results, i.e. the abstract state at a program location, to check for difference (we temporarily ignore the question of program location matching for the sake of argument). As we are dealing with over-approximations, it would be impossible to claim equivalence based on two separate states: say we only wish to analyze patch (2), and would like to compare the two states right after entering the loop. Optimally, these states will both contain the $i > 0$ and $i < 0$ constraints (we will later

```

1 void foo(int arr[], unsigned len) {
2     int i = len;
3     i--;
4     while (i) {
5         arr[i] = i;
6         i--;
7     }
8 }

1 void foo(int arr[], unsigned len) {
2     int i = len;
3     if (len == 0) return; (1)
4     i--;
5     while (i) {
6         arr[i] = i--; (2)
7     }
8 }

```

Figure 1. Example looping code and patched version

describe in detail the handling of non-convex data). Though it would be tempting to claim that the programs are equivalent at this point, it would be wrong, as each program may have arrived at the constraints by entirely different means. For instance, the $i < 0$ constraint in one state could be the result of a (convexly) joining $-5 \leq i < 0$ and $i \leq -10$, but it could be the result of joining $i = 0$ and $i \leq -1$ in the other. This key observation, which basically states that *equality under abstraction does not assure concrete equality*, dictates the use of a combined analysis as otherwise we can never hope to establish equivalence.

One option for performing this combined, correlating analysis is by defining a special correlating semantics which requires performing a dual analysis on both programs, whilst maintaining abstract information regarding the two sets of variables together. This is a viable option, however we chose a different approach where we would simply construct a single correlating program, denoted $P \bowtie P'$ (for the correlation of a program P and its patched version P'), that will hold the variables and statements of the two programs. These will be interleaved in such a way where matching statements (that appear in both versions) will be adjacent, thus allowing the analysis to maintain equivalence. We keep the variables of the programs separated by tagging all variables from the newer version. We opted for the correlating program solution since it allows us to employ standard analysis frameworks [1] and abstract domains [5]. Another advantage is that the correlating program building process supplies us with a matching of program locations, thus we are able to check for difference at appropriate locations. Lastly, since $P \bowtie P'$ is a syntactically correct program, that contains the semantic of both programs, we are able to use existing techniques of symbolic execution and equivalence checking, as used in previous work [4, 9?], **to achieve potentially better results**, as we allow for a more fine-grained checking and differencing (as opposed to input-output checking). We will elaborate on these issues, and on the process of building the correlating program in Section ??.

Having defined a facility for performing a joint analysis, we need to define our abstraction in such a way where we are able to maintain equivalence (under abstraction) if such exists, or provide a precise description of the difference while maintaining soundness. Considering a standard relational abstraction, the only case where equivalence in variables is assured, is when both versions of the variable equal the same concrete value. As this is usually not the case (especially for unknown inputs) we explicitly force the abstraction to initially assume equivalence until proven otherwise. For example, when analyzing patch (2) from Fig. 2.1 (ignoring patch (1)) while forcing initial equivalence $i = i'$ as when we enter the loop (i' is the new version of i), the analysis will advance over the

first lines of the loops and will temporarily lose equivalence but keep the $i = i' + 1$ constraint which will be used to restore equivalence as it moves past the second loop line in the first program. This is a key feature of our analysis as we allow *temporary equivalence divergence* with the ability to later restore it.

As stated, explicit equivalence over abstraction is key for proving equivalence. Consequently, we cannot use a standard convex abstract domain [5], we will quickly lose equivalence and will not be able to restore it. Let us examine the two pieces of equivalent code in Fig. 2.1. The optimal result here would be that $\llbracket P \bowtie P' \rrbracket^\# \models_{\equiv_V}$, but the use of a convex domain here can never yield this result.

```

1  if (x>3) y = 4;

1  z = 4;
2  if (x>3) y = z;

```

Figure 2. Example conditional code and patched version

3. Preliminaries

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

Program A program P, \dots . We are able to convert any program to this format... We exclude recursion for now.

Program Location A program location $loc \in Loc$,

Program Label A program label $lab \in Lab$, is a unique identifier for a certain location in a program. Every location has a label (usually the program counter). We also define two special labels for the start and exit locations of the program as *begin* and *fin* respectively.

Concrete State A concrete program state is a tuple $\sigma \equiv \langle loc, values \rangle \in \Sigma$ mapping the set of (integer) program variables to their concrete (integer) 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 $\llbracket P \rrbracket$.

Concrete Trace A program trace $\pi \in \Sigma^*$, is a sequence of states $\sigma_0, \sigma_1, \dots$ 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 $\Pi \subseteq \Sigma^*$. We refer to these semantics as concrete state semantics. We also define the following standard operations on traces:

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

4. Semantics of Union Program

After defining the notion of concrete state and trace delta, we require a facility for producing said delta (alongside a variable and label correspondence) from a pair of programs P, P' . Delta can be produced by either (i) separately analyzing the programs and computing delta postmortem, between abstract states computed at matching labels (ii) analyze programs *together* obtaining a joint correlating abstract state and extracting delta from it. In this section we will explore method (ii) for producing our delta. **We will later on expand on option (ii) advantages over option (i) which encouraged us to choose it.**

4.1 General Product Program

A simple approach for a joint analysis is to construct a product program $P \times P'$ where at every point during the execution we can perform a program step (as defined in Definition 3) of either programs. The product program has a duo-state (σ, σ') and each step updates (σ, σ') accordingly. The product program can also be seen as a concurrent run $P || P'$ where every interleaving is possible. The product program emphasizes the fact that, as described in Section 2, the notion of Δ is unclear without an established variable and label correspondence. Choosing the location where Δ is checked is a key part of identifying differences. Consider Fig. 4.1, which presents a product automata of the simple program with itself, we see that even in this trivial program, although it is clear that $\Delta = \emptyset$, checking for difference in any of the non-correlating states will result in a false difference being reported. As this example demonstrates, selecting a correct label correspondence is crucial for a meaningful delta, we will elaborate on our approach for choosing DP in Subsection 4.2.

```

1  void foo() {
2      int x = 0;
3  }

```

Figure 3. Program P

4.2 Program Correspondence and Differential Points

Selecting the point where Δ is computed is vital for precision. As mentioned, a natural selection for diff points would be at the endpoints of traces but that loses meaning under the collecting semantics. A possible translation of this notion under the collecting semantics would be to compute delta between *all* the endpoints of the two programs i.e. $DP = \{(fin, fin') | fin \in exit(P), fin' \in exit(P')\}$ somehow differentiating the final states of the programs. This approach is problematic for two reasons:

1. Comparing all endpoints results in a highly imprecise delta. This is shown by the simple exercise of taking program with 2 endpoints and comparing it with itself.
2. This choice for DP may result in missing key differences between versions. If at some point during the calculation existed a delta that failed reaching the final state - it will be ignored. An interesting example for this is an array index receiving different bounds after a patch (but later overwritten so that it is not propagated to some final state).

Alternatively, the brute force approach where we might attempt to capture more potential diffs by selecting a diff-point after every line, will result in a highly inaccurate result as, for instance in Fig. 4.1, many diffs will be reported although there is no difference. Finally, we must be careful with the selection of DP as it affects the soundness of our analysis: we might miss differences if we did not correctly place diff-points in locations where delta exists. **Our approach employs standard syntactic diff algorithm ?? for producing the correlation. This selection for DP assures soundness.** The Diff approach works well since two versions of the same software (and especially those that originate from subsequent check-ins to a code repository) are usually similar. Another important factor in the success of the diff is the guarded instruction format for our programs (as defined in Definition 3). Transforming both programs to our format helps remove a lot of the "noise" that a patch might introduce yet it is superior to low lever intermediate representation as it retains many qualities (such as variable names, conditions, no temporaries, etc.). **See Appendix ?? for examples illustrating said benefits and qualities.** There are alternate ways for creating the correspondence such as **graph equivalence, etc.**,

this could be a subject of future research. Calculating delta according to DP over the product automata is a complex task as it allows both programs to advance independently. We formulate the *union program* as a restricted product automata where we advance the programs while keeping the correlation allowing for a superior calculation of delta using our correlating abstract domain later defined in Section 6.

4.3 The Union Program $P \cup P'$

We will generally describe the process of constructing the union program. A more elaborate and formal description of the algorithm can be found in Algorithm ?? . The union program is an optimized structure where not all pairs of (σ, σ') are considered, but only pairs that result from a controlled execution, where correlating instructions (according to DP) in P and P' will execute together. This will allow for superior precision. As said, the main idea is to create one program which contains both versions. The union starts out as (exactly) the older version P (after being converted to our guarded instruction form). Afterwards a syntactic diff with P' (also transformed to guarded mode) is computed (the programs are not combined just yet). In fact, this is the point where DP is created as the diff supplies us with the correlation between labels we desire. Then P' 's instructions are interleaved into the guarded P while maintaining the correlation found by the diff (matched instructions will appear sequentially). Just before a patched instruction is interleaved into the union, all variables that appear in it are tagged, as to make sure that the patched instructions will only affect patched variables. Thus we maintain the semantics of running both programs correctly while achieving a new construct that will allow us to analyze change more easily and precisely. **Fig. ?? holds a complete union program of the program in Fig. ?? and its patched version, a graphic description as a controlled automata is shown in Fig. ??** . Note that if we view the general union program as a concurrent program, then this optimized program can be viewed as a partial-order reduction applied over the concurrent program. One final observation regarding the union program is that it is a legitimate program that can be run to achieve the effect of running both versions. This ability allows us to use dynamic analysis and testing techniques such as fuzzing ?? and directed automated testing ?? which may produce input that lead to states approximated by Δ .

4.4 Analyzing Union Programs

Analysis of a guarded union program has certain caveats. In order

```

1  1: guard g = (i>0);
2    if (g) i--;
3    if (g) goto 1;

```

Figure 4. example program illustrating guard analysis caveat

to correctly analyze the program in Fig. 4.4 we need our analysis to assume $(i > 0)$ whenever taking the true branch on the `if (g)` instruction and $(i \leq 0)$ when taking the false branch. However, since the `i--` instruction invalidates this assumption we would need to update the guard assumption to $(i > -1)$ which would complicate the analysis as we would need to consider updating the guard assumption while widening etc. Our solution simply incorporates the guard's assumption the first time it encounters the guard and allows it to flow to the rest of the nodes. We are not in danger of losing the assumption during the following join as our join employs a partition-by-equivalence strategy and will not join the two states where $g, i > 0$ and $\neg g, i \leq 0$.

5. Concrete Semantics

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

5.1 Concrete State Differencing

Comparing two different programs P and P' under our concrete semantics means comparing *traces*. A trace is composed of concrete states thus we first need a way to compare states. Since a state σ is a mapping $Var \rightarrow Val$ we need a correspondence between variables in P (i.e. Var) and those in P' (i.e. Var').

Variable Correspondence A variable correspondence $VC \in Var \leftrightarrow Var'$, is a partial mapping between 2 sets of program variables. Any variable in Var may be matched with one in Var' and vice versa. Naturally, the comparison between states will occur between values mapped to corresponding variables, as described by VC .

Our experience suggests that in most cases the set of variables stays the same over subsequent versions ($Var = Var'$) and in cases where Var does change, it's by the addition of new variables or removal of an old one. Therefore we concluded that for our purposes, matching variables by name is sufficient for finding a precise difference. Thus we define our standard correlation to be: $VC_{Eq} \equiv \{v \mapsto v' | v \in Var \wedge name(v) = name(v')\}$ and vice versa. Take notice that in cases where a variable was removed (added) by a patch, it will still be correlated to its patched (unpatched) version, although it does not exist, and we will report it as part of the difference - this will adhere to our definition of difference. In cases where the matching is not that straightforward, for example, when variables name change, the correspondence can be generated in [a more sophisticated way](#) .

Concrete State Delta Given two concrete states $\sigma \in \Sigma_P, \sigma' \in \Sigma_{P'}$ and a variable correspondence VC , we define the concrete state delta $\Delta_S(\sigma, \sigma')$ of σ and σ' as the part of the state σ where corresponding variables do not agree on values (with respect to σ'). Formally: $\Delta_S(\sigma, \sigma') \equiv \{(var, val) | var \in Var \wedge VC(var) = var' \wedge \sigma(var) = val \neq \sigma'(var')\}$.

We note that the state delta is not necessarily symmetric. In fact, the direction in which it is used has meaning in the context of a program P and a patched version of it P' . We define $\Delta_S^- = \Delta_S(\sigma, \sigma')$ which means the part of the state that was "changed" or "removed" in P' and $\Delta_S^+ = \Delta_S(\sigma', \sigma)$ which stands for the part added in P' .

For instance, given two states $\sigma : \{x = 1, y = 2, z = 3\}$ and $\sigma' : \{x' = 0, y' = 2, w' = 4\}$ and assuming our default VC then $\Delta_S^- = \{x = 1, z = 3\}$ since x and x' match and do not agree on value, y and y' agree (thus are not in delta) and no data exists for z' in σ' . Respectively, $\Delta_S^+ = \{x' = 0, w' = 4\}$. We believe that Δ_S is a precise and natural definition for difference as it indeed captures the change made in the patched version. We will later on extend the definition to describe difference between abstract states.

We defined a notion for difference between states but this is insufficient to describe difference between whole runs of programs i.e. traces. Naturally, we are only interested in traces that originate from **the same input** (for the input variables correlated by VC) although one can contemplate the question of delta between *any* two traces of a program (or two programs) however this is not the problem addressed in this paper and every mention of trace differentiation will assume the traces agree on input. The way to differentiate traces is by differentiating their states, but which states? this is not a trivial question since traces can vary in length and order of states. We need a state correspondence for choosing the states to be differentiated within the two traces. We define it as following:

Trace Diff Points Given two traces π and π' , we define a trace index correspondence relation named trace diff points denoted DP_π as a matching of indexes specifying states in where concrete state delta should be computed. Formally: $DP_\pi \equiv \{(i, i') | i \in 0..|\pi|, i' \in 0..|\pi'|\}$.

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

Trace Delta Given traces π, π' of programs P, P' respectively, and a state correspondence DP_π we define the trace delta $\Delta_T(\pi, \pi')$ as state differentiations between all corresponding states in π and π' . Formally, for every $(i, i') \in DP_\pi$, Δ_T will contain the mapping $i \mapsto \Delta_S(\sigma_i, \sigma'_{i'})$, thus the result will map certain states in π to their state delta with the corresponding state' in π' (deemed interesting by DP_π). We note that Δ_T^+ and Δ_T^- have a similar meaning except that now it's in a trace context.

One possible choice for DP_π would be the endpoints of the two traces $\{(fin, fin')\}$ (assuming they are finite) meaning differentiating the final states of the executions or formally: $\Delta_{fin}^- \equiv \Delta_{Fin}(\pi, \pi') = \{fin \mapsto \Delta_S(\sigma_{fin}, \sigma'_{fin'})\}$ (we will also be interested in Δ_{Fin}^+). It is clear that the delta is not sufficient for truly describing the difference between said traces as it will only compare final state values and will miss out on what happened during the execution. This can be overcome by instrumenting the semantics such that the state will contain all "temporary" values for variable along with the location the values existed (by, for instance, adding temporary variables), making all the differences visible at the end point. **As a design choice we chose not to do this (since it creates many variables, choosing VC here is hard, etc)**.

Defining the diff-points over any two traces is a daunting task since two traces of two separate (although similar) programs can vastly differ. If we take a look at two versions of a program depicted in Fig. 5.1 and the following traces generated from the input $x = 2$: $\pi = \{x = 2, i = 0\}\{i = 1\}\{i = 2\}$ and $\pi' = \{x' = 2, i' = 0\}\{i' = 1\}\{i' = 2\}\{i' = 3\}\{i' = 4\}$ (we omit labels and only mention parts of the trace where variable values change), we see that even in this simple program, finding a correlation based on traces alone is hard. However, one can get the sense that using program location as a means of correlation can produce a meaningful result. For example, if we look at all the possible values for i in label lab and differentiate them from the values in the patched version (in the same location), we get the meaningful result that i in the patched version can range from $x + 1$ up to $2x$. The need to find a label-based match, along with the need for abstraction to allow scalability, is filled by moving past the concrete trace semantics to a concrete *collecting semantics*.

<pre> 1 void foo(unsigned x) { 2 unsigned i = 0; 3 lab:if (i >= x) return; 4 i++; 5 goto lab; 6 } </pre>	<pre> 1 void foo(unsigned x) { 2 unsigned i = 0; 3 lab:if (i >= 2*x) return; 4 i++; 5 goto lab; 6 } </pre>
---	---

Figure 5. P, P' differentiation candidates

5.2 Differencing at Program Labels

Collecting Semantics of a Program Label Given a label $l \in Lab$ in a program P we define the collecting semantics of a program label $states(l) \subseteq \llbracket P \rrbracket$ as all the concrete states that are possible at that label (i.e. exist in some trace reaching that label). Formally:

1. $at \equiv Lab \rightarrow (\Pi \rightarrow 2^{\llbracket P \rrbracket^*})$ for a given label l and trace π , is the set of prefixes of π that end in a state labeled l formally: $at(l, \pi) \equiv \{\pi * \pi' \in pre(\pi) \wedge last(\pi') = l\}$.

2. $states(l) \equiv \{last(at(l, \pi)) | \pi \in \Pi\}$.

Now our problem is reduced to matching the collecting semantics of labels $states(l)$. Again we are encountered with the question of how to match these labels. We note that the trace indexing correspondence DP_π defined in Definition 5.1 is no longer useful here as we need to differentiate sets of states belonging to a certain program label $states(l)$. Thus we require a correspondence or *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: $DP_{Lab} \equiv \{(l, l') | l \in Lab, l' \in Lab'\}$.

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 selecting DP_{Lab} in a correct and meaningful way in Subsection 4.2. We will briefly mention that in order to find a meaningful correspondence, we employ a standard syntactic diff algorithm [?] on the two versions of the program. Our experiments show that for the purpose of differentiating versions of program, this approach works well since patched versions tend to be syntactically similar (especially when the versions come from two successive check-ins to the code repository).

Given the correspondence of labels in P and P' (from DP), we now know which $states(l)$ and $states(l')$ we need apply delta on. Delta is now applied on *sets* of states, which is defines by simply applying the state delta Definition 5.1 between each of the states in both states and removing from the result matched states or formally: $\Delta_C(states(l), states(l')) \equiv \{\sigma \in states(l) | \neg \exists \sigma' \in states(l') \cdot \Delta_S(\sigma, \sigma') = \emptyset\}$

For example, given two sets of states $C : \{\sigma_1 : \{x = 0, y = 0\}, \sigma_2 : \{x = 1, y = 2\}\}$ and $C' : \{\sigma'_1 : \{x = 0, y = 0\}, \sigma'_2 : \{x = 4, y = 5\}\}$ and using VC_{Eq} then $\Delta_C^- = \{\{x = 1, y = 2\}\}$ and $\Delta_C^+ = \{\{x = 4, y = 5\}\}$. Note that Δ_C^+ now obtains the meaning of "lost states" as in states which existed in the previous version and removed by the patch (similarly Δ_C^- here means "new states").

We must remember however, that the sets of states to be compared are potentially unbounded which means that the delta we compute may be unbounded too. Therefore we must use an abstraction over the collecting semantics that will allow us to represent the collecting semantics in a bounded way.

6. Abstract Correlating Semantics

In this section, we introduce our correlating abstract domain which allows bounded representation of union program state while focusing on maintaining equivalence between correlated variables. This comes at the cost of an acceptable lose of precision of other numerical information of the variables. We represent variable information using standard numerical abstract domain. To allow for temporary divergence of equivalence (due to union program structure) we keep a set of abstracts (as divergence is non-convex). As we will show, this allows for restoration of equivalence later on (if indeed equivalence holds) and in the case we are unable to converge, we will record the precise state information (which will produce a more precise diff) before aggressively joining the abstracts set into one abstract, continuing the analysis and avoiding exponential blow-up. We start off by abstracting the collecting semantics in Subsection ??.

In the following, we assume an abstract numerical domain $ND = \langle NC, \sqsubseteq_{ND} \rangle$ equipped with operations \sqcap_{ND} and \sqcup_{ND} , where NC is a set of numerical constraints over the variables in Var , and do not go into further details about the particular abstract domain. We also assume that the numerical domain ND allows for

a sound over-approximation of the concrete collecting semantics (given a sound interpretation of program operations).

Correlating Abstract State A correlating abstract program state σ^h , is a tuple $\langle l, G^{nc} \rangle \in \Sigma^h$, where G^{nc} is a **group** of numerical constraints, each capturing relationships between numerical variables of both the original and patched programs P and P' . The semantics of $G^{nc} = nc_1, nc_2, \dots, nc_n$ is $nc_1 \wedge nc_2 \wedge \dots \wedge nc_n$ where each nc_i is a disjunction of numerical constraints. Let us explicitly define the operations of the domain:

- $G_1 \sqsubseteq_{CD} G_2 \iff \forall nc_1 \in G_1 \exists nc_2 \in G_2 : nc_1 \sqsubseteq_{ND} nc_2$
- $G_1 \sqcap_{CD} G_2 \equiv nc_1 \sqcap_{ND} nc_2 | nc_1 \in G_1 \wedge nc_2 \in G_2$
- $G_1 \sqcup_{CD} G_2 \equiv G_1 \cup G_2$

One major advantage of our correlating domain over using two separate domains, is the ability to preserve equivalence in the face of non-linear operations - this argument may be too thin to include.

For example, if we take our motivating example from Fig. ?? and annotate it with our correlating domain, we will get after line XXX the state XXX and after line XXX the state XXX. This example emphasizes the need for our correlating domain to hold a group of numerical constraints since Although precise, the sets of abstracts produced by interpreting our example using our correlating domain are not very informative. One cannot easily deduce whether equivalence is kept or if the state holds a difference. To answer this question we define the *correlating abstract state delta*.

6.1 Correlating Abstract State Differencing

Given a state in our correlating domain, we want to compute the version difference, if exists, at that state. However, since our input now is one correlating state holding information of both versions of variables, there is no straightforward way of defining the difference (unlike previous delta definitions Definition ??, Definition ??). We overcame this by treating the numerical constraints in our domains as geometrical objects and formulating delta based on that.

Correlating Abstract State Delta Given two abstract states and a correspondence VC , the correlating state delta $\Delta_A(\sigma^h, \sigma'^h)$, computes abstract state differentiations between σ^h and σ'^h . The result is an abstract state $\sqsubseteq \sigma^h$ approximating all concrete values possible in P but not in P' (regarding variables that match in VC). Formally, the delta is simply $\sigma^h \setminus \sigma'^h$ but since this concept is vague to the reader and furthermore, does not exist in most domain implementation (and specifically in the ones we used) we break it down to a simpler multi-step operation as following:

1. $U \equiv \sigma^h \sqcap \sigma'^h$ is the joint state of the original and patched program. No precision is lost while joining states as they operate on different variables. This state, as well as the ability to cleanly separate variables, is achieved by symbolically executing a *union program* as defined in Section 4.
2. R is a state abstracting the concrete states shared by the original and patched program. It is achieved by computing: $R \equiv U|_{V=V'} \equiv U \sqcap \bigwedge \{v = v' | VC(v) = v'\}$.
3. $R|_V$ is the projected state where all the variables from Var' are eliminated from R .
4. $\bar{R}|_V$ is the negated state i.e. $D \setminus R|_V$ and it is computed by negating $R|_V$ (as mentioned before, all logical operations, including negation, are defined on our representation of an abstract state).
5. Eventually: $\Delta(\sigma^h, \sigma'^h) \equiv \sigma^h \sqcap \bar{R}|_V$ meaning it is part of the original program state σ^h that does not appear in σ'^h i.e. appears

in the negation of R (which is the intersection of both abstract states).

A geometrical representation of Δ_A calculation can be seen in Fig. ??

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

6.2 Minimizing Divergence

One can see from our motivating example that it is not feasible to allow our correlating domain to keep diverging and double in size with every conditional as it will exponentially blow up the analysis run-time and memory. Instead, we employ an equivalence conserving canonization technique such that after every join will either (i) check for equivalence in the joined state and if it is kept, join all of G^{nc} 's abstracts into one abstract, potentially losing precision but preserving equivalence. (ii) see that equivalence is not kept and allow it to converge, for now. In order for use to truly maintain equivalence after option (1) is executed, our domain must not lose precision of variable equality over join. Option (ii) entails that once equivalence is broken, our analysis will quickly explode in memory as we will no longer be able to minimize sets of constraints by joining all. We solve this by performing our canonicalization at the next diff-point, but make sure to record the "exploded" state before joining all abstracts as it is potentially *more precise* and may be used to produce a more informative delta (without losing soundness).

7. Evaluation

We mainly tested our tool on the GNU core utilities, differencing versions 6.10 and 6.11. This benchmark included 40 patches where most of the patches (35) were a one-line patch aimed at updating the version information string in the code. Our analysis easily showed equivalence for these programs. About 10 of these patches included actual changes to numerical variables and we were 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 union compiler named *ucc* which creates union programs from any two C programs as well as a differencing oriented dataflow analysis solver for analyzing union programs, both tools use the LLVM and Clang compiler infrastructure. We analyze C code directly thus benefiting from a low number of variables as there are no temporary values as there might appear in an intermediate representation. We also benefit from our delta being computed over original variables. As mentioned in Section 4, we normalize the input programs before unifying them for a simpler analysis.

Analysis of some of our benchmarks required the use of widening. We applied a basic widening strategy which widens all cfg blocks once reaching a certain threshold. All of our experiments were conducted running on a Intel(R) Core-i7(TM) processor with 4GB.

7.1 Results

Tab. 1 summarizes the results of our analysis. The columns indicate the benchmark name and description, lines of code for the analyzed program, the number of lines added and removed by the patch, the number of diff-points generated, the numerical domain used, and the number of differences found in the analysis (i.e. number of diff-points where $\Delta \neq \emptyset$). In our benchmarks, we focused on computing intra-procedural difference between the two versions of

Table 1. Experimental Results				
Name	Domain	#Added	#Deleted	#DiffPoints
dd.i	ppl	52	54	87
id.i	ppl	15	6	26
pr.i	ppl	10	3	13
su.i	ppl	2	2	2
env.i	ppl	2	2	3
seq.i	ppl	10	10	15
nice.i	ppl	3	3	36
test.i	ppl	3	3	14
chmod.i	ppl	3	3	7
nohup.i	ppl	2	2	18
paste.i	ppl	24	17	4
rmdir.i	ppl	3	0	3
users.i	ppl	3	4	30
chroot.i	ppl	3	3	22
md5sum.i	ppl	15	7	32
runcon.i	ppl	2	2	3
lbracket.i	ppl	3	3	14
setuidgid.i	ppl	7	4	34
chown-core.i	ppl	3	3	10

procedures. Procedure calls presented difficulty as they potentially change global variables and local variables through pointers. We overcame this by either (i) assuming equivalence (alone) once we encounter a call to a procedure we already established as equivalent or (ii) warn that all results regarding variables touched by the procedure is un-sound. **In the majority of our benchmarks we identified calls only to library and system procedures thus we could omit their effect as they do not change variables beyond those given as parameter or those being assigned the return value.** All differences reported describe, in constraints over variables, an existing delta at that program point.

Identifying delta down the line One advantage of our analysis method is the ability to identify differences in variables that were not directly affected by the patch. Fig. 7.1 shows part of the `bsd_split_3` function that was patched by the line marked by a comment. Note that although the patch directly restricts the value range of `s_len` to above zero, our analysis is able to identify the effects on the index variable `i` and also report a lost program state of `i > -1` further down the line...

Non-convex delta

Maintaining equivalence in loops

8. Related Work

Bounded symbolic execution in CLang As prior work we used the CLang infrastructure [1] static analysis graph reachability engine in order to perform a simple and bounded state differentiation exploration. We used the existing infrastructure and it's abstract representation facilities to simply record every location where the 2 versions of the variables differ. This of course was not sufficient since it only presents a bounded solution and we will show the limitations of this method by example.

Existing work on patch-based exploit generation Brumley, Poosankam, Song and Zheng [?] is the prominent work addressing patch-based analysis. We differ from this work in the following aspects:

1. First, the problem definition in said work is different from our own. They aim to find an *exploit* for vulnerabilities fixed by a certain patch. Furthermore, this exploit is defined in relevance to a *securitypolicy* which can differ. While our goals are similar to those of [?], we achieve them by solving 2 extended problems of a) recording the delta between new variable values

```

1 static bool
2 bsd_split_3(char *s, size_t s_len, unsigned char **hex_digest, char **
3 size_t i;
4
5 if (s_len == 0) return false; // Patch
6
7 *file_name = s;
8
9 /* Find end of filename. The BSD 'md5' and 'sha1' commands do not es
10 filenames, so search backwards for the last ')'. */
11
12 while (i < s_len - 1;
13 while (i && s[i] != ')')
14 i--;
15
16 if (s[i] != ')') return false;
17
18 s[i++] = '\0';
19
20 while (ISWHITE(s[i])) i++;
21
22 if (s[i] != '=') return false;
23 i++;
24
25 while (ISWHITE(s[i])) i++;
26
27 *hex_digest = (unsigned char *) &s[i];
28 return true;
29 }

```

Figure 6. md5sum.c bsd_split_3 function patch

and old ones and b) producing input from said values. These problems are a superset of the problem described in [?] and solving them has the potential for a much more complete and sound result.

2. We aim to find differentiation between every variable changed by the patch and analyze that differentiation while they concentrate on input sanitation alone. Thus if in the patched program some variable has changed in a way that does not involve input validation, it will be disregarded: for instance if an array index variable i to a buffer B is patched by adding an assignment $i = \text{sizeof}(B) - 1$, it will be ignored in the previous work while we will record that the old version of i can no longer have values greater than $\text{sizeof}(B) - 1$ and may use it for exploit generation.
3. We perform our analysis on the source code of the program and patch instead of the binary. Working on a higher level gives us much more data thus potentially allowing for more results.

Kroening and Heelan [?] main focus focus was producing an exploit from given input that is known to trigger a bug. No patch is involved in the process. Our goal is to produce said input from the corrected software thus [?] can be used to create an exploit from our results.

Song, Zhang and Sun [?] also relate to the patch-based exploit generation problem but their main focus is on finding similarities between versions of the binary to better couple functions from the original program with their patched counter-part, a problem that was not addressed in [?]. Also their method of recognizing possible exploits is degenerate and relies on identifying known input validation functions that were added to a certain path - a method that could be easily overcome.

Oh [?] presented a new version for the DarunGrim binary diffing tool aimed at better reviewing patched programs and specifically finding patches with security implications. The goal of the tool is to help researches who manually scan patches for the purpose of producing intrusion prevention system signatures. The tool relies mainly on syntactic analysis of patterns to produce a security

implication score for procedures patches making them a candidate for manual inspection. [?] used the DarunGrim binary diffing tool EBDs for their experiment.

Person, Dwyer, Elbaum and Pasareanu [8] introduced an extension and application of symbolic execution techniques that computes a precise behavioral characterization of a program change called differential symbolic execution. As we also implemented bounded symbolic execution as our preliminary work we will discuss this method in comparison to our own.

Godlin and Strichman [?] developed a method for proving the equivalence of similar C programs under certain restrictions based on and existing functional verification tool. This was a basis for future work regarding equivalence and we intend to base our work upon these advances.

Kawaguchi, Lahiri and Rebelo [7] defined the concept of *conditional equivalence* meaning under which conditions (inputs) are 2 different versions of a program equivalent (i.e. produce the same output). Their goal is to keep software changes from breaking procedure contracts and changing module behavior too drastically and they achieve this by computing the conditions under which the behavior is preserved. This work indirectly addresses our problem and we believe we can leverage their techniques for producing the inputs that break the equivalence while focusing on bug triggering rather than contract breaking.

[?]

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. Appendix

A.1 Algorithm 2 : Convert P To Guarded Instruction Format

The algorithm is constructive i.e. it takes a procedure P and outputs the new lines for a guarded version of P . The original P is not part of the output.

- *Stage 0:* Output P 's signature.
- *Stage 1:* Convert all `while` constructs to `if` and `goto` constructs.
- *Stage 2:* For each non branch instruction I :
If I is a declaration, output it under a new block.
Otherwise collect all branch conditions C under which I executes. Produce the code line: `if ($\bigwedge C$) I;`

We take a small but sufficient example to demonstrate the algorithm:

```

1 void foo(char Out[], int n) {
2     char In[42];
3     int i;
4     if (n >= 42) n = 41;
5     while (i < n) {
6         int j = i + 1;
7         In[i] = Out[j];
8         i++;
9     }
10 }
```

Stage 1 Output:

```

1 void foo(char Out[], int n) {
2     char In[42];
3     int i;
4     if (n >= 42) n = 41;
5 1: if (i < n) {
6         int j = i + 1;
7         In[i] = Out[j];
8         i++;
9         goto 1;
10    }
11 }
```

Stage 2 Output:

```

1 void foo(char Out[], int n) {
2     char In[42];
3     int i;
4     if (n >= 42) n = 41;
5     {
6         int j;
7 1:     if (i < n) j = i + 1;
8         if (i < n) In[i] = Out[j];
9         if (i < n) i++;
10        if (i < n) goto 1;
11    }
12 }
```

Figure 7. Algorithm 1 application example

Note that as required, the loop is implemented in means of branches and goto and that the branches do not exceed the 1 level of nesting allowed.

A.1.1 Algorithm 2 Assumptions

1. We only deal with `while` loop statements.
2. Logical statement (branch conditions) do not have side-effects i.e. they do not change variable values.

A.2 Algorithm 1 Example

We run algorithm 1 on the previous example with the following patch to line 4:

```

- if (n >= 42) n = 41;
+ if (n >= 42) n = 40;
```

We start off from the output of algorithm 2 i.e. a "guarded" program that performs simple buffer handling:


```

1 void foo(char Out[], int n) {
2     char In[42];
3     int i;
4     if (n >= 42) n = 41;
5
6     {
7         int j;
8     l:   if (i < n) j = i + 1;
9         if (i < n) In[i] = Out[j];
10        if (i < n) i++;
11        if (i < n) goto l;
12    }
13
14 }

```

The result:

```

1 void foo(char Out[], int n, char Out'[], int n') {
2     char In[42], In'[42];
3     int i, i';
4
5     if (n >= 42) n = 41;
6     if (n' >= 42) n' = 40;
7
8     {
9         int j, j';
10    l:   if (i < n) j = i + 1;
11    l':  if (i' < n') j' = j' + 1;
12        if (i < n) In[i] = Out[j];
13        if (i' < n') In'[i'] = Out'[j'];
14        if (i < n) i++;
15        if (i' < n') i'++;
16        if (i < n) goto l;
17        if (i' < n') goto l';
18    }
19 }

```

Figure 8. Algorithm 1 Application Example

B. Worklist

B.1 Points to Hammer

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

B.2 TODO

1. show me an example of where comparing abstract values at the end is not sound (e.g., we had some example with intervals that demonstrates this).

B.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' ?

References

- [1] clang: a c language family frontend for llvm, 2007.
- [2] 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–.
- [3] CHEN, L., MINÉ, A., WANG, J., AND COUSOT, P. Interval polyhedra: An abstract domain to infer interval linear relationships. In

Proceedings of the 16th International Symposium on Static Analysis (Berlin, Heidelberg, 2009), SAS '09, Springer-Verlag, pp. 309–325.

- [4] GODLIN, B., AND STRICHMAN, O. Regression verification. In *DAC* (2009), pp. 466–471.
- [5] JEANNET, B., AND MINÉ, A. Apron: A library of numerical abstract domains for static analysis. In *CAV* (2009), pp. 661–667.
- [6] KAWAGUCHI, M., LAHIRI, S. K., AND REBELO, H. Conditional equivalence. Tech. rep., MSR, 2010.
- [7] LAHIRI, S., HAWBLITZEL, C., KAWAGUCHI, M., AND REBÊLO, H. Symdiff: A language-agnostic semantic diff tool for imperative programs. In *CAV* (2012).
- [8] PERSON, S., DWYER, M. B., ELBAUM, S. G., AND PASAREANU, C. S. Differential symbolic execution. In *SIGSOFT FSE* (2008), pp. 226–237.
- [9] RAMOS, D., AND ENGLER, D. Practical, low-effort equivalence verification of real code. In *Computer Aided Verification*, vol. 6806 of *LNCS*. Springer, 2011, pp. 669–685.
- [10] 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.