

# An Extensible Verified Validator For Simple Optimizations in LLVM

Sungkeun Cho\*, Joonwon Choi\*, Jeehoon Kang\*,  
Chung-Kil Hur, and Kwangkeun Yi

{skcho,jwchoi}@ropas.snu.ac.kr, {jeehoon.kang,gil.hur}@sf.snu.ac.kr, kwang@ropas.snu.ac.kr  
Seoul National University

September 30, 2014

## Abstract

Recent advances in theorem proving technology have made it possible to write a completely bug-free compiler such as CompCert. However, this technology have not been fully applied to mainstream compilers such as GCC and LLVM due to the huge amount of work required. Nevertheless, we believe that the verified validation technique will reduce verification efforts significantly and thus make it feasible to provide a formal guarantee of correctness for the full LLVM compiler.

As a first step towards the grand goal, we have developed an extensible verified validator for simple optimizations in LLVM, mainly targeted for all micro-optimizations (about 400) in the instruction combine pass. Our validator, based on a simple form of relational Hoare reasoning, can be used to validate various kinds of simple optimizations without losing completeness (*i.e.*, it succeeds to validate all correct results at least by design). We achieve both generality and completeness, not by automation, but by extensibility, which allows to freely add inference rules that are proven to be correct.

## 1 Introduction

Since compilers are indispensable tools for software development, their reliability is crucial, especially for safety-critical software. However, as recent researches [11, 3] have shown, mainstream C compilers such as GCC and LLVM contain plenty of bugs that matter practically. For example, [11] found 79 bugs in GCC and 202 bugs in LLVM, and [3] found 79 bugs in GCC and 68 bugs in LLVM.

In recent years, there has been a considerable amount of work to increase the reliability of a compiler. The main techniques include random testing [11, 3], translation validation [6, 7], verified compilation [4, 8] and verified validation [8, 9, 10]. These approaches have their own advantages and disadvantages. The former two are general-purpose techniques, so that they have been used to find bugs or validate compilation results in mainstream compilers. However, these methods do not provide formal guarantee of absence of bugs. On the other hand, the latter two techniques provide the highest level of guarantee, a formal proof of semantics preservation. However, they have not been fully applied to mainstream C compilers such as GCC and LLVM, due to the huge amount of work required.

The grand goal of our project, called LLVMBERRY, is to provide a formal guarantee of semantics preservation for the LLVM compiler, by minimizing the amount of verification effort required. To this end, we will use the verified validation technique only without using any verified compilation, and try to reduce verification efforts by maximizing the merits of the validation approach. As the first step towards our grand goal and as a case study, we have

\*The first three authors equally contributed to this paper and are listed alphabetically.

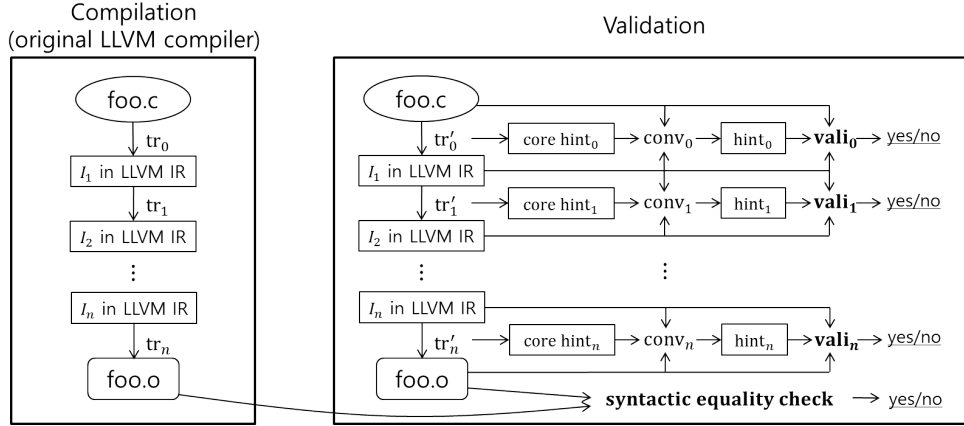


Figure 1: The framework of LLVMBERRY: Programs in boldface need to be verified.

developed a verified validator, which we call SIMPLEBERRY, for numerous micro optimizations in the instruction combine pass of the LLVM compiler, which is the main result of the paper. We elaborate more on these points in the following.

## 1.1 Verified Validation vs. Verified Compilation

Before we discuss why we chose verified validation over verified compilation, we briefly review the two methods.

Compiler verification consists of two steps: (i) to implement a compiler as a computable function  $C : \text{Src} \rightarrow \text{option Tgt}$  in a proof assistant such as Coq, and (ii) to prove that every target program produced by the compiler implements only the behaviors specified by the source program:

$$\forall s, t. C(s) = \text{Some } t \implies \text{behaviors}(t) \subseteq \text{behaviors}(s) .$$

Then, one can extract the verified compiler function as an executable program and use it to compile source programs.

Verified validation for a particular compiler, introduced by Tristan and Leroy [8, 9, 10], also consists of two steps: (i) to implement a translation validator specialized for the compiler as a computable function  $V : \text{Src} \times \text{Tgt} \rightarrow \text{bool}$  in a proof assistant, and (ii) to prove that any translation validated by the validator preserves the semantics of the source program:

$$\forall s, t. V(s, t) = \text{true} \implies \text{behaviors}(t) \subseteq \text{behaviors}(s) .$$

Then, one can extract the verified validator as an executable program and use it to check whether a compilation result of the compiler is valid or not.

Verified validation for the LLVM compiler has the following advantages over verified compilation.

1. We do not need to re-implement the compiler in a proof assistant. This is a big advantage because fully implementing the LLVM compiler in a proof assistant like Coq is practically impossible due to the huge amount work required. Even if it is possible, since the verified compiler is extracted into a functional language like OCaml, the verified compiler is hard to be as fast as the original compiler written in C++. This may be a problem because the compilation speed matters in practice.
2. A compiler development team and a verification team can be separated. More specifically, in the verified validation approach, compiler experts can write a compiler in C++ and

verification experts can write a verified validator in a proof assistant, whereas, in the verified compilation approach, those with expertise in both compilation and verification have to write a verified compiler in a proof assistant. This separation is important in practice because not many people have expertise in both area.

3. Verifying a validator requires less verification efforts than verifying a compiler. This is because, in many cases, a validator that checks translations of a particular optimization pass is much simpler than the optimization pass itself. Furthermore, in some cases, a single validator can be reused to validate results of other similar optimization passes.
4. Verified validation is less sensitive to changes to the compiler. Since verified validation does not directly verify the code of a compiler, it may not be much affected by changes to the compiler, such as version upgrades, which regularly occur in mainstream compilers. More specifically, even if an optimization algorithm changes in order to improve the performance, the validation algorithm may remain the same or change slightly.

Now we discuss the downsides of verified validation and how they can be overcome.

1. A verified validator may fail to validate correct translations. According to the aforementioned correctness statement, verification of a validator only guarantees that the validator never validates any incorrect translation. Thus, there is no guarantee that the validator validates all correct translations.

However, since we develop a validator that is targeted for a particular optimization pass, it is always possible to design the validator to be *complete w.r.t. the target optimizer*, meaning that it validates all correct results that are intended by the optimizer. Although the design may have a bug that leads to incompleteness, such a bug is easily detected by the verified validator. More specifically, if the validator fails to validate a translation result of the target optimizer, there are only two possibilities: (i) the translation is incorrect (*i.e.*, the optimizer has a bug); or (ii) the validator is incomplete (*i.e.*, the design has a bug). In either case, the details of the validation will guide us to the point where the bug occurs, so that we can easily fix the bug.

2. Validation may take long. Since validation includes checking various invariants and conditions, it may be much slower than compilation. If validation is invoked for every compilation, this will increase the whole compilation time, which can be problematic.

However, compilation and validation can be completely separated. Thus, a developer can first compile a source program and use the compiled program before the validation of the compilation is completed. Moreover, one can selectively validate compilation results during software development. For example, she can validate a compilation result only when the compiled program behaves erroneously, or when it is a release version.

Also, validation time can be reduced greatly by using multi-cores or distributed computing since validation is highly parallelizable. A compiler consists of many optimization passes, each of which may be performed many times during compilation. Thus, a single compilation consists of at least hundreds of optimization steps. Since each optimization step can be validated completely independently, validating the whole optimization steps is perfectly suited to multi core and distributed computing environment.

## 1.2 Our Approach to Verified Validation for LLVM

The framework of LLVMBERRY is depicted in Figure 1. In this framework, we provide a completely separated validation module (RHS of the figure), so that a user can use the original LLVM compiler for compilation (LHS of the figure) and *a posteriori* validate the compilation result whenever she wants.

In the validation module,  $tr'_i$  is the original optimization pass  $tr_i$  of the LLVM compiler, extended with additional code that generates *core hints* for validation. Core hints are the key information about why the optimizer performed such optimization, and they are used by the associated validator  $val_i$  to determine validity of the optimization result. It is such hints that make it possible to design a *complete* validator for any particular optimizer.

$conv_i$ , written in OCaml, is a *hint converter* that converts core hints in JSON format to data that the validator  $val_i$  can understand. Furthermore,  $conv_i$  can also expand core hints to full hints by adding more hints inferable from the core hints together with the input and the output programs. Note that  $conv_i$  is not a part of the validator  $val_i$  because  $conv_i$  does not need to be verified.

$val_i$  is a *validator* that takes input and output programs with hints and determines validity of the translation using the hints. Unlike other parts,  $val_i$  has to be verified in order for our validation to be sound. Thus,  $val_i$  is first written and verified in the proof assistant Coq, and then extracted as an executable program in OCaml. The formal soundness statement for  $val_i$  is as follows:

$$\forall s, t, h. \text{val}_i(s, t, h) = \text{true} \implies \text{behaviors}(t) \subseteq \text{behaviors}(s).$$

The final part of the validation is to check whether the original output program ('foo.o' in LHS of the figure) and the validated output program ('foo.o' in RHS of the figure) are syntactically equal. Rigorously speaking, this syntactic equality checker also has to be verified, but we just use the GNU "diff" because we trust it for such a simple task.

Our basic principle for minimizing verification efforts is not to care about validation time. In other words, whenever we can reduce verification efforts by increasing validation time, we do it. We think this principle is reasonable because, as we discussed before, validation can be performed off-line and is highly parallelizable.

The following are examples of our specific efforts that would make verification easier, though they may increase validation time.

1. We try to split a translation into as small translations as possible. This would make verification simpler because a validator just needs to focus on a small translation at a time.
2. We try to use the same validator for many different optimizations by designing it as general as possible.
3. We do not prove the preservation of properties such as well-formedness and read-onlyness of a function that can be additional verification overhead for each validator. Instead, we just dynamically check the properties at validation time whenever they are needed. Though we have to verify the correctness of the well-formedness and read-onlyness checker, they are verified once and used many times.

### 1.3 Extensible Verified Validator for Simple Optimizations

To support our arguments for verified validation of the LLVM compiler, we developed an extensible verified validator, called SIMPLEBERRY, and used it to validate micro optimizations in the instruction combine pass of the LLVM version 3.0. For formalization of the LLVM IR (Intermediate Representation), we used the Vellvm [12, 13] framework.

Here are summaries of our work.

- SIMPLEBERRY is based on a simple form of relational Hoare reasoning and designed to be general enough to validate many simple optimizations that do not change the control flow graph (CFG) of programs, including all micro optimizations (about 400) in the instruction combine pass.

- We did not perform any complex automation to achieve such generality because such automation would require complex verification. Instead, we designed SIMPLEBERRY to be extensible, so that one can freely add user-defined inference rules specific to a particular optimization by just proving the correctness of the rules.
- Examples of optimization that we validated do not only include peephole optimizations, but also simple inter-block, heap and phi-node optimizations.
- As of now, we have validated 50 optimizations out of about 400 optimizations in the instruction combine pass. However, after we built the infrastructure for SIMPLEBERRY and tested it for 3 representative optimizations, it took 5 person-weeks to add 47 more optimizations. Adding one micro optimization includes understanding the optimization code, adding and verifying new inference rules, adding code that generates core hints, and sometimes debugging. We expect to add the remaining 350 optimizations within 4 months by two PhD students.

The contributions of the paper are as follows:

1. We designed a framework for verified validation of LLVM compiler.
2. We applied the verified validation approach to a mainstream compiler for the first time.
3. We developed an extensible verified validator for the first time.

## 2 Overview

In this section, we present the high-level design of SIMPLEBERRY.

Since SIMPLEBERRY does not essentially rely on the details of LLVM IR such as the SSA property, we will present our work using a simple abstract language, where the syntax and the semantics are intuitively clear. First, we assume that our language has variables (also called registers)  $x, y, z, \dots$  and the heap. We use C-like notation for heap access such as  $*x = y$  and  $x = *y$ . We will explain other details of the language on the fly when we need them.

SIMPLEBERRY is basically a proof checker for a simple form of relational Hoare logic. For validation, we first give an invariant at each line of code and ask SIMPLEBERRY to check whether these invariants are valid (or consistent). Since checking validity of a Hoare logic is undecidable in the presence of heap, we allow a user to add any sound inference rules in order to increase validation power of SIMPLEBERRY.

**Simple Relational Invariants** First of all, SIMPLEBERRY has a fixed heap invariant. Throughout the reasoning, informally speaking, we implicitly keep the following invariant for the heaps  $h_{\text{src}}$  and  $h_{\text{tgt}}$  in the source and the target programs:

$$\begin{aligned} \exists h_{\text{src}}^e, h_{\text{src}}^i, h_{\text{tgt}}^e, h_{\text{tgt}}^i. \\ (h_{\text{src}}^e \uplus h_{\text{src}}^i \subseteq h_{\text{src}}) \wedge (h_{\text{tgt}}^e \uplus h_{\text{tgt}}^i \subseteq h_{\text{tgt}}) \wedge (h_{\text{src}}^e \approx h_{\text{tgt}}^e). \end{aligned} \quad (1)$$

This says that the source and the target heaps are split into two parts such that one parts are equivalent and the other parts are unrelated. We call the former *public heaps* and the latter *isolated heaps*. (See Section 5 for more formal details.)

SIMPLEBERRY allows a user to specify an invariant about variables at each line of code (*i.e.*, pre and post conditions of Hoare logic). An invariant at line  $i$  have three components: source property  $P_i$ , target property  $Q_i$  and may-diff set  $D_i$ .

$$i : \begin{array}{ccc} \{ P_{i-1} \} & \{ Q_{i-1} \} & \{ D_{i-1} \} \\ \text{src command}_i \mapsto \text{tgt command}_i & & \\ \{ P_i \} & \{ Q_i \} & \{ D_i \} \end{array}$$

The source property  $P_i$  is a property about the state of the source program that should hold just after the execution of the instruction at line  $i$ . Similarly,  $Q_i$  is a property that holds at line  $i$  in the target program.  $D_i$  is what connects the states of the source and the target programs at line  $i$ .

Source and target properties  $P_i$  and  $Q_i$  are sets consisting of five kinds of basic properties: (i) equality about a register value  $x = e$  (e.g.,  $x = y + 1$ ); (ii) equality about a heap value  $*p = v$ ; (iii) no alias between pointer values  $p \neq q$ ; (iv) allocated pointers  $\text{alc}(p)$ ; and (v) isolated pointers  $\text{isol}(p)$ . The semantics of (i) and (ii) is as expected.  $p \neq q$  denotes that the values contained in the variables (or registers)  $p$  and  $q$  are pointer values that point to disjoint parts of the heap.  $\text{alc}(p)$  denotes that the variable  $p$  contains a pointer value that points to an allocated part of the heap.  $\text{isol}(p)$  denotes that  $x$  contains a pointer value that points to the isolated heap.

A may-diff set  $D_i$  at line  $i$  contains all variables  $x$  such that the value of  $x$  in the source program and the value of  $x$  in the target program may be different at line  $i$ . In other words, every variable that is not in  $D_i$  should contain equivalent values in the source and the target programs.

Regarding the notion of equivalent values and heaps, we can simply understand them as the same values and heaps for now. Additionally, if equivalent values are pointer values, they should point to the public heaps in the source and the target. Similarly, the values in the public heap  $h_{\text{src}}^e$  (resp.,  $h_{\text{tgt}}^e$ ) should satisfy that if they are pointer values, they point back to the public heap  $h_{\text{src}}^e$  (resp.,  $h_{\text{tgt}}^e$ ). More formal details will be presented in Section 5.

Finally, we allow logical (or existential) variables in invariants in order to increase the expressive power. A logical variable  $\hat{x}$  has nothing to do with the physical value of the variable  $x$  and thus is distinguished from the variable  $x$ . It only denotes the existence of some value and is used to connect several equations. We will explain more details in Section 3.3.

**Validation Algorithm** Given invariants, SIMPLEBERRY checks the consistency of the invariants. At each line, it first conservatively checks whether the source and the target programs raise the same side effects, such as heap update, function call and error. Then, it calculates a post condition  $P'_{i-1}, Q'_{i-1}, D'_{i-1}$  that is strong enough for our purpose<sup>1</sup>, and tries to check whether  $P'_{i-1}, Q'_{i-1}, D'_{i-1}$  implies  $P_i, Q_i, D_i$ .

$$\begin{array}{ccc}
 \{ P_{i-1} \} & \{ Q_{i-1} \} & \{ D_{i-1} \} \\
 i : \text{src command}_i \mapsto \text{tgt command}_i & & \\
 \left[ \begin{array}{c} P'_{i-1} \\ P''_{i-1} \end{array} \right] & \left[ \begin{array}{c} Q'_{i-1} \\ Q''_{i-1} \end{array} \right] & \left[ \begin{array}{c} D'_{i-1} \\ D''_{i-1} \end{array} \right] \\
 \{ P_i \} & \{ Q_i \} & \{ D_i \}
 \end{array}$$

Checking the implication is the most difficult task. Here, SIMPLEBERRY does not perform any automation but, instead, simply applies the inference rules given as hints to derive a new invariant  $P''_{i-1}, Q''_{i-1}, D''_{i-1}$ . Finally, it checks whether  $P''_{i-1} \supseteq P_i$ ,  $Q''_{i-1} \supseteq Q_i$  and  $D''_{i-1} \subseteq D_i$  hold. Note the direction of containment is reversed for may-diff sets because they are properties about those variables not in the may-diff sets.

**Extensibility** In order to validate a particular optimization, one may need special inference rules for that optimization. Since it is infeasible for SIMPLEBERRY to initially have all such inference rules, we allow to register user-defined inference rules. In this way, we can easily extend the power of SIMPLEBERRY to validate results of a new optimizer.

We do not have any restriction on the expressibility of user-defined rules. An inference rule can be any computable function, say  $f$ , in Coq that takes an invariant  $(P, Q, D)$  and returns

<sup>1</sup>It is not strongest because we do not perform any complex automation.

**Transformation**

B1
9: $x = a + 1 \mapsto x = a + 1$
B4
5: $y = x + 2 \mapsto y = a + 3$

**Invariant Commands**

9@B1: propagate\_to 5@B4 in  
src

**Inference Rule Commands**

5@B4: add\_assoc  $x y$  in src  
5@B4: remove\_maydiff  $y$

**Validation**

B1	$\{ \}$	$\{ \}$	$\{ \}$
9: $\mathbf{x = a + 1}$	$\mapsto$	$\mathbf{x = a + 1}$	
$[x = a + 1]$		$[x = a + 1]$	$[ ]$
$\{x = a + 1\}$		$\{ \}$	$\{ \}$
B4	$\{x = a + 1\}$		$\{ \}$ $\{ \}$
5: $\mathbf{y = x + 2}$		$\mapsto$	$\mathbf{y = a + 3}$
$[x = a + 1, y = x + 2]$			$[y = a + 3]$ $[y]$
add_assoc $x y a 1 2 3$ in src			
$[x = a + 1, y = x + 2, y = a + 3]$			$[y = a + 3]$ $[y]$
remove_maydiff $y (a + 3)$			
$[x = a + 1, y = x + 2, y = a + 3]$			$[y = a + 3]$ $[ ]$
$\{ \}$			$\{ \}$ $\{ \}$

Figure 2: Validation of Add Assoc Optimization.

a new invariant  $(P', Q', D')$ . SIMPLEBERRY simply applies this inference rule according to the hint command. The correctness of the rule  $f$  is the following:

$$\forall S : \text{State}. \forall P, Q, D. S \models (P, Q, D) \implies S \models f(P, Q, D).$$

The correctness theorem of SIMPLEBERRY requires the correctness of inference rules as a precondition. Thus, one has to prove them in Coq in order to have the end-to-end formal guarantee. However, in case she is very confident about the correctness of an inference rule, she can also put the correctness into the trust base by making it as an axiom.

**Overall Process** The overall process of validation, according to the LLVMBERRY framework (see Figure 1), is as follows. First, the modified compiler  $tr'$  performs an optimization and generates the core hints that consists of (i) commands for generating invariants, and (ii) commands for generating inference rules. Then, from the core hints, the hint converter  $conv$  generates invariants and inference rules as full hints, and feed them to the validator  $vali$ . Finally, the validator checks the validity of the given translation.

### 3 Validation Examples

In this section, we walk you through the details of SIMPLEBERRY with various examples.

#### 3.1 Add Assoc Optimization

Our first example is the add assoc optimization. Consider the example optimization in Figure 2, where  $x = a + 1$  at line 9 of block B1 followed by  $y = x + 2$  at line 5 of block B4 is optimized to  $x = a + 1$  followed by  $y = a + 3$ .

For the validation, we first need to add the following two rules, where  $FV(e)$  is the set of all variables occurring in  $e$ . The correctness of the rules is intuitively clear.

```

add_assoc  $x\ y\ a\ c_1\ c_2\ c_3$  in src :=  $\lambda(P, Q, D)$ .
  if  $(x = a + c_1 \in P) \wedge (y = x + c_2 \in P) \wedge c_3 = c_1 + c_2$ 
  then  $(P \cup \{y = a + c_3\}, Q, D)$  else  $(P, Q, D)$ 
remove_maydiff  $x\ e$  :=  $\lambda(P, Q, D)$ .
  if  $x = e \in P \wedge x = e \in Q \wedge FV(e) \cap D = \emptyset$ 
  then  $(P, Q, D \setminus \{x\})$  else  $(P, Q, D)$ 

```

Second, the modified compiler generates the core hints, consisting of invariant and inference rule commands, in Figure 2. The hint converter then generates invariants and inference rules according to the commands. The invariant command `9@B1:propagate_to 5@B4 in src` inserts the equation  $x = a + 1$  at line 9 of B1 into the source property of every line between line 9 of B1 and line 5 of B4 (see the black curly bracket in Validation of Figure 2). The command `5@B4:add_assoc  $x\ y$  in src` infers the remaining arguments and inserts the inference rule `add_assoc  $x\ y\ a\ 1\ 2\ 3$  in src` at line 5 of B4, and then `5@B4:remove_maydiff  $y$`  inserts `remove_maydiff  $y\ (a + 3)$`  at the same line.

Finally, the validator checks the validity of the invariants. At each line, it first generates a strong post condition. For example, it generates  $[x = a + 1, y = x + 2]$ ,  $[y = a + 3]$ ,  $[y]$  at line 5 of B4. It adds the invariant  $y = x + 2$  in the source and  $y = a + 3$  in the target because of the assignments executed. Also, it adds  $y$  to the may-diff set because it cannot simply determine whether the newly assigned values to  $y$  in the source and the target are the same. Unlike at line 5 of B4, it does not add  $x$  to the may-diff set in line 9 of B1 because the simple equality checker of SIMPLEBERRY can infer that the values in  $x$  in the source and the target are the same. Indeed, the equality checker simply checks whether the two instructions are syntactically the same and the free variables in the RHS are not in the may-diff set.

At line 9 of B1, the validator trivially checks that the strong post condition implies the invariant at the line. At line 5 of B4, it sequentially applies the `add_assoc` and `remove_maydiff` rules according to the given hints and gets the final invariant  $[x = a + 1, y = x + 2, y = a + 3]$ ,  $[y = a + 3]$ ,  $[\ ]$ . This invariant trivially implies  $\{ \}, \{ \}, \{ \}$ . Thus, we checked the consistency of all invariants, which completes the validation.

## 3.2 Heap Optimizations

The next examples are heap optimizations: store-load optimization and dead store and alloca elimination.

**Store-Load Optimization** For the store-load optimization in Figure 3, we add to SIMPLEBERRY the following rule, whose correctness is easy to show.

```

trans_via_heap  $x\ y\ p$  in src :=  $\lambda(P, Q, D)$ .
  if  $*p = x \in P \wedge *p = y \in P$ 
  then  $(P \cup \{x = y\}, Q, D)$  else  $(P, Q, D)$ 

```

Then the modified compiler generates the core hints in Figure 3. The command `5@B3:propagate_to 7@B3 in src` adds the property `alc( $p$ )` between 5@B3 and 7@B3 in the source; the command `6@B3:propagate_to 9@B3 in src` adds `* $p$  =  $y$`  between 6@B3 and 9@B3 in the source; and `7@B3:propagate_neq_to  $q\ p$  8@B3 in src` adds  `$q \neq p$`  between 7@B3 and 8@B3 in the source. The command `9@B3:trans_via_heap  $x\ y$  in src` and `9@B3:remove_maydiff  $x$`  add the rules `trans_via_heap  $x\ y\ p$`  and `remove_maydiff  $x\ y$`  at line 9 of B3 in the source.

The validator now checks the validity of the invariants. At line 5, SIMPLEBERRY adds `alc( $p$ )` in the post condition, but does not add  $p$  to the may-diff set because it can simply infer



that  $p$  contains equivalent pointers in the source and the target that point to allocated blocks of the heaps. At line 6, SIMPLEBERRY infers that  $*p = y$  holds. At line 7, it infers, in addition to  $\text{alc}(q)$ , that  $q \neq p$  holds (*i.e.*,  $q$  and  $p$  point to disjoint parts of the heap) because  $\text{alc}(p)$  holds in the pre condition and  $q$  is newly allocated.

At line 8, due to the heap update, SIMPLEBERRY does several things. First, it checks that the variables  $q$  and  $z$  are not in the may-diff set in order to make sure that the global heap invariant is preserved. Second, it does not remove  $*p = y$  despite of heap update because there is no alias between  $p$  and  $q$  due to  $q \neq p$ . Finally, it adds  $*q = z$  in the source and the target due to the heap update performed.

At line 9, SIMPLEBERRY adds  $*p = x$  and  $x = y$  in the source and the target, and then adds  $x$  to the may-diff set because it is not easy to see whether the values assigned to  $x$  are equivalent. Then, the following inference rules properly update the invariants and completes the validation, as shown in Figure 3. Here, it is important to note that since we dropped the condition  $\text{alc}(p)$ , we do not know that  $p$  is allocated and thus the instruction  $x = *p$  may fail. However, it is ok because if the source program produces an error, the target program is allowed to behave arbitrarily.

**Dead Store And Alloca Elimination** Thanks to the store-load optimization, line 6 and 5 became dead code and thus can be eliminated by dead store elimination followed by dead alloca elimination. This is an example that uses the isolated set invariant. For presentation purpose, we perform the two eliminations at the same time (see Figure 4).

First of all, note that we insert instructions called *logical nop* in the target program. These are a part of the hint and used only for validation purpose. Using logical nop, we can align the source and the target program.

The modified compiler generates the core hints in Figure 4. The command `*:propagate_maydiff_to p *` adds  $p$  to the may-diff set at every line, and `5@B3:propagate_isol_to p 6@B3` in `src` adds  $\text{isol}(p)$  at line 5 of B3.

SIMPLEBERRY then checks the validity of the invariants. At line 5, it treats the allocated heap in the source as an isolated heap because there is no corresponding allocation in the target, and thus adds  $\text{isol}(p)$  in addition to  $\text{alc}(p)$ . At line 6, SIMPLEBERRY can make sure that the heap update in the source does not break the global heap invariant because the heap update is made to the isolated heap due to  $\text{isol}(p)$ . The rest is straightforward.

### 3.3 Fold-Phi-Bin Optimization

The final example is the fold-phi-bin optimization, which optimizes a phi-node and is sensitive to the control flow of the program. In this example, we will see how logical variables are used.

A phi-node is an assignment that assigns different values depending on the incoming block. For example, in the program of Figure 5, the block B3 has two incoming blocks, B2 and B3. The instruction  $c = \phi(\mathbf{a@B2}, \mathbf{b@B3})$  assigns to  $c$  the value of  $a$  if coming from B2, and that of  $b$  if coming from B3.

The fold-phi-bin transformation, shown in Figure 5, defers the increment-by-one computation after the phi node, using a temporary variable  $t$ . In order to align the source and the target programs, we insert a logical nop in the source side. Since the block B3 behaves differently depending on the incoming block, SIMPLEBERRY validates B3 separately for each incoming block. In Figure 5,  $\text{B3} \leftarrow \text{B2}$  is the validation of B3 when coming from B2, and  $\text{B3} \leftarrow \text{B3}$  is the validation when coming from B3.

First, we need to add the following inference rules.

```

trans_via_var x y e in src :=  $\lambda(P, Q, D)$ .
  if  $x = y \in P \wedge y = e \in P$ 
  then  $(P \cup \{x = e\}, Q, D)$  else  $(P, Q, D)$ 

```

```

replace_var  $x\ y\ z\ e\ e'$  in  $\text{src} := \lambda(P, Q, D)$ .
  if  $x = y \in P \wedge z = e \in P \wedge e' = e[y/x]$ 
  then  $(P \cup \{z = e'\}, Q, D)$  else  $(P, Q, D)$ 

intro_eq  $x\ y := \lambda(P, Q, D)$ .
  if  $x \notin D \wedge \hat{y} \notin \text{Fv}(P, Q, D)$ 
  then  $(P \cup \{x = \hat{y}\}, Q \cup \{x = \hat{y}\}, D)$  else  $(P, Q, D)$ 

```

Here, an interesting rule is `intro_eq`, which allows to introduce redundant equalities using a new fresh logical variable. Note that  $\hat{y} \notin \text{Fv}(P, Q, D)$  means that the logical variable  $\hat{y}$  does not occur in  $P$ ,  $Q$  and  $D$ , *i.e.*, it is fresh. Here the newly introduced equalities say that there is some value represented by  $\hat{y}$  that is equal to the value of  $x$  in the source and the target, and furthermore the values of  $\hat{y}$  in the source and the target are equivalent because  $\hat{y} \notin D$ . We will see below how such redundant equalities are used in this example.

For validation, the modified compiler generates hints that produce the invariants and inference rules shown in Figure 5. We omit how they are generated since it is similar as in the previous examples. Note that there is a single post condition of block B3 that merges the post conditions of B3←B2 and B3←B3.

The validation of B2 is straightforward. For validation of B3, SIMPLEBERRY first checks that the post condition of the block B2 implies the pre condition of block B3←B2, and the post condition of B3 implies the pre condition of B3←B3. Also it checks that the post conditions of B3←B2 and B3←B3 both imply the common post condition of block B3. In validation of B3←B2, there is nothing special except that SIMPLEBERRY adds  $c = a$  and  $t = x$  after execution of  $\mathbf{c} = \phi(\mathbf{a}@B2, \mathbf{b}@B3)$  and  $\mathbf{t} = \phi(\mathbf{x}@B2, \mathbf{c}@B3)$  because we assume that the incoming block is B2.

In the validation of B3←B3, at line 0 SIMPLEBERRY adds  $c = b$  and  $t = c$ . However, in the source, we cannot simply add  $c = b$  because the pre condition already contains equalities about  $c$ , whose value is not the same as the newly assigned value of  $c$ . Here, SIMPLEBERRY uses a logical variable  $\hat{c}$  to represent the old value of the variable  $c$ . It simply renames the variable  $c$  in the pre condition to  $\hat{c}$ .<sup>2</sup> Since we assigned a value to  $c$  only in the source, SIMPLEBERRY has to add  $c$  and  $\hat{c}$  to the may-diff set. However, as a side effect, we lose the information that the value of  $\hat{c}$  in the source is equivalent to that of  $c$  in the target. It is the redundant equations  $\hat{c} = \hat{u}$  in the source and  $c = \hat{u}$  in the target that recovers the lost information since  $\hat{u}$  is not in the may-diff set. At line 1, after renaming  $c$  to  $\hat{c}$  and adding  $c = t + 1$  in the target, SIMPLEBERRY regains, via the logical variable  $\hat{u}$ , the fact that  $c$  contains the equivalent value in the source and the target. Reasoning at line 2 is straightforward. Thus, the validation succeeds.

## 4 Semantics

Before we explain the algorithm and the verification, we first present a formal semantics. For the presentation purpose, we abstract away the unnecessary details of LLVM IR. In other words, we present a semantics for an idealized LLVM IR.

### 4.1 Program Semantics

$C = (G, F) \in \text{Config}$  is a static configuration initialized after loading a program  $\mathcal{P}$ , where  $G \in \text{Id} \rightarrow \text{option Addr}$  is a global variable map that maps a variable name to the associated heap address, and  $F \in \text{Addr} \rightarrow \text{option Fdef}$  is a function lookup table that maps a function pointer to its declaration and definition. Here, `option`  $T$  is the usual option type consisting of either  $\perp$  or  $[t]$  for some  $t \in T$ . A function name is also a global variable, so that you can

<sup>2</sup>Here, before renaming  $c$  to  $\hat{c}$ , SIMPLEBERRY first removes every predicate containing  $\hat{c}$  from the pre condition. Though there will not be any such predicate because the program is in SSA, this removal makes the verification of SIMPLEBERRY easier.

get the definition of a function, say `foo`, by first looking up `foo` in  $G$  (i.e.,  $[p] = G(\text{foo})$ ) and then get the definition via  $F$  (i.e.,  $[\text{fd}] = F(p)$ ). Note that the static configuration  $\mathcal{C}$  remains the same during the execution. Throughout this section, by  $\mathcal{C}_{\text{src}}$  and  $\mathcal{C}_{\text{tgt}}$  we denote the static configurations of the source and the target program.

$\mathcal{S} = (\sigma, h) \in \text{State}$  is a *program state* that changes over time during the execution, and consists of a *call stack*  $\sigma \in \Sigma (= \text{list } \mathbb{F})$  and a *heap*  $h \in \text{Heap} (= \text{Addr} \rightarrow \text{Value} \uplus \{\text{dead}\})$ . Here, a call stack  $\sigma$  is a list of stack frames and a heap  $h$  is a partial function that maps heap addresses to either values or dead cells (dead). An allocated cell becomes dead when it is freed, and is never used for allocation again. A *stack frame*  $(pc, f, s) \in \mathbb{F}$  consists of a program counter  $pc$ , the current function  $f$ , and the local store  $s \in \text{Store} (= \text{Id} \rightarrow \text{Value})$ . A program counter  $pc (= n @ B_i)$ , consisting of a block label  $B_i$  and a line number  $n$ , indicates the current code pointer being executed inside  $f$ . The local store  $s$  maps variables to values. A value  $v \in \text{Value} (= \text{Int} \uplus \text{Addr})$  is either an integer, or an address (pointer value).

The small-step operational semantics is a judgment of the form:

$$\mathcal{C} \vdash \mathcal{S} \rightarrow^\varepsilon \mathcal{S}' ,$$

which denotes that in the static configuration  $\mathcal{C}$ , the current state  $\mathcal{S}$  steps to a next state  $\mathcal{S}'$ , emitting an observable IO event  $\varepsilon$ . You can understand this as a usual small-step semantics.

As we have seen in the previous examples, we need to align instructions using logical nops. We can easily derive an instrumented semantics that takes into account a hint  $N$  for the locations of logical nops, which has the following form:

$$(\mathcal{C}, N) \vdash \mathcal{S} \rightarrow^\varepsilon \mathcal{S}' .$$

## 4.2 Hint Semantics

As we have seen in Section 2, we have two kinds of invariants: the implicit heap invariant and invariants at each line of code.

**Implicit Heap Invariant** As discussed in Section 2, the implicit heap invariant relating the source heap  $h_{\text{src}}$  and the target heap  $h_{\text{tgt}}$  partitions them into public and isolated heaps ( $h_{\text{src}}^e \uplus h_{\text{src}}^i \subseteq h_{\text{src}}$  and  $h_{\text{tgt}}^e \uplus h_{\text{tgt}}^i \subseteq h_{\text{tgt}}$ ). To represent such an invariant we need partitioning information  $\alpha^3$ ,  $\beta_{\text{src}}$ , and  $\beta_{\text{tgt}}$ , where  $\alpha$  is a bijection between the addresses of  $h_{\text{src}}^e$  and  $h_{\text{tgt}}^e$ ,  $\beta_{\text{src}}$  is the set of addresses in  $h_{\text{src}}^i$ , and similarly for  $\beta_{\text{tgt}}$ . Of course,  $\beta_{\text{src}}$  should be disjoint from  $\alpha$ , and the same for  $\beta_{\text{tgt}}$ .

In addition to the partitioning, we need more properties on the heaps.

- Every cell in the isolated heaps should not be dead:

$$\begin{aligned} & (\forall a \in \beta_{\text{src}}. h_{\text{src}}(a) \neq \text{dead}) \wedge \\ & (\forall a \in \beta_{\text{tgt}}. h_{\text{tgt}}(a) \neq \text{dead}) . \end{aligned}$$

- The values in the corresponding locations of the public heaps w.r.t.  $\alpha$  should be equivalent modulo  $\alpha$ , denoted  $h_{\text{src}}^e \stackrel{\alpha}{\simeq} h_{\text{tgt}}^e$ , if:  $\forall (a_{\text{src}}, a_{\text{tgt}}) \in \alpha$ .

$$(h_{\text{src}}^e(a_{\text{src}}) = h_{\text{tgt}}^e(a_{\text{tgt}}) = \text{dead}) \vee (h_{\text{src}}^e(a_{\text{src}}) \stackrel{\alpha}{\simeq} h_{\text{tgt}}^e(a_{\text{tgt}})) ,$$

where values  $v_{\text{src}}$  and  $v_{\text{tgt}}$  are equivalent modulo  $\alpha$ , denoted  $v_{\text{src}} \stackrel{\alpha}{\simeq} v_{\text{tgt}}$ , if:

$$(v_{\text{src}} = v_{\text{tgt}} \in \text{Int}) \vee ((v_{\text{src}}, v_{\text{tgt}}) \in \alpha) .$$

We write  $(h_{\text{src}}, h_{\text{tgt}}) \in \mathcal{H}(\alpha, \beta_{\text{src}}, \beta_{\text{tgt}})$  when the heap invariant holds for  $h_{\text{src}}$  and  $h_{\text{tgt}}$ .

<sup>3</sup>The notion of memory injection  $\alpha$  is originally used in the verification of CompCert [5].

**Transformation**

<b>B3</b>	
5: $\mathbf{p} = \mathbf{alloca}(1)$	$\mapsto \mathbf{p} = \mathbf{alloca}(1)$
6: $\ast\mathbf{p} = \mathbf{y}$	$\mapsto \ast\mathbf{p} = \mathbf{y}$
7: $\mathbf{q} = \mathbf{alloca}(1)$	$\mapsto \mathbf{q} = \mathbf{alloca}(1)$
8: $\ast\mathbf{q} = \mathbf{z}$	$\mapsto \ast\mathbf{q} = \mathbf{z}$
9: $\mathbf{x} = \ast\mathbf{p}$	$\mapsto \mathbf{x} = \mathbf{y}$

**Invariant Commands**

5@B3: propagate\_to  
7@B3 in src  
6@B3: propagate\_to  
9@B3 in src  
7@B3: propagate\_neq\_to  
 $q\ p\ 8@B3$  in src

**Inference Rule Commands**

9@B3: trans\_via\_heap  $x$   
 $y$  in src  
9@B3: remove\_maydiff  $x$

**Validation**

<b>B3</b>			
	{ }	{ }	{ }
5: $\mathbf{p} = \mathbf{alloca}(1)$	$\mapsto \mathbf{p} = \mathbf{alloca}(1)$		
	[ alc( $p$ ) ]	[ alc( $p$ ) ]	[ ]
	{ alc( $p$ ) }	{ }	{ }
6: $\ast\mathbf{p} = \mathbf{y}$	$\mapsto \ast\mathbf{p} = \mathbf{y}$		
	[ alc( $p$ ), $\ast p = y$ ]	[ $\ast p = y$ ]	[ ]
	{ alc( $p$ ), $\ast p = y$ }	{ }	{ }
7: $\mathbf{q} = \mathbf{alloca}(1)$	$\mapsto \mathbf{q} = \mathbf{alloca}(1)$		
	[ alc( $p$ ), $\ast p = y$ , alc( $q$ ), $q \neq p$ ]	[ alc( $q$ ) ]	[ ]
	{ $\ast p = y$ , $q \neq p$ }	{ }	{ }
8: $\ast\mathbf{q} = \mathbf{z}$	$\mapsto \ast\mathbf{q} = \mathbf{z}$		
	[ $\ast p = y$ , $q \neq p$ , $\ast q = z$ ]	[ $\ast q = z$ ]	[ ]
	{ $\ast p = y$ }	{ }	{ }
9: $\mathbf{x} = \ast\mathbf{p}$	$\mapsto \mathbf{x} = \mathbf{y}$		
	[ $\ast p = y$ , $\ast p = x$ ]	[ $x = y$ ]	[ $x$ ]
	trans_via_heap $x\ y\ p$ in src		
	[ $\ast p = y$ , $\ast p = x$ , $x = y$ ]	[ $x = y$ ]	[ $x$ ]
	remove_maydiff $x\ y$		
	[ $\ast p = y$ , $\ast p = x$ , $x = y$ ]	[ $x = y$ ]	[ ]
	{ }	{ }	{ }

Figure 3: Validation of Store-Load Optimization.

## Transformation

<b>B3</b>
5: $\mathbf{p} = \mathbf{alloca}(1) \mapsto \mathbf{logicalnop}$
6: $\mathbf{*p} = \mathbf{y} \mapsto \mathbf{logicalnop}$
7: $\mathbf{q} = \mathbf{alloca}(1) \mapsto \mathbf{q} = \mathbf{alloca}(1)$
8: $\mathbf{*q} = \mathbf{z} \mapsto \mathbf{*q} = \mathbf{z}$
9: $\mathbf{x} = \mathbf{y} \mapsto \mathbf{x} = \mathbf{y}$

## Invariant Commands

$\mathbf{*} : \mathbf{propagate\_maydiff\_to}$   
 $\mathbf{p} \ \mathbf{*}$   
 $\mathbf{5@B3} : \mathbf{propagate\_isol\_to}$   
 $\mathbf{p} \ \mathbf{6@B3}$  in src

## Inference Rule Commands

## Validation

<b>B3</b>			
	{ }	{ }	{ $p$ }
5: $\mathbf{p} = \mathbf{alloca}(1) \mapsto \mathbf{logical\ nop}$			
	[ $\mathbf{alc}(p), \mathbf{isol}(p)$ ]	[ ]	[ $p$ ]
	{ $\mathbf{isol}(p)$ }	{ }	{ $p$ }
6: $\mathbf{*p} = \mathbf{y} \mapsto \mathbf{logical\ nop}$			
	[ $\mathbf{isol}(p), \mathbf{*p} = \mathbf{y}$ ]	[ ]	[ $p$ ]
	{ }	{ }	{ $p$ }
7: $\mathbf{q} = \mathbf{alloca}(1) \mapsto \mathbf{q} = \mathbf{alloca}(1)$			
	[ $\mathbf{alc}(q)$ ]	[ $\mathbf{alc}(q)$ ]	[ $p$ ]
	{ }	{ }	{ $p$ }
8: $\mathbf{*q} = \mathbf{z} \mapsto \mathbf{*q} = \mathbf{z}$			
	[ $\mathbf{*q} = \mathbf{z}$ ]	[ $\mathbf{*q} = \mathbf{z}$ ]	[ $p$ ]
	{ }	{ }	{ $p$ }
9: $\mathbf{x} = \mathbf{y} \mapsto \mathbf{x} = \mathbf{y}$			
	[ $\mathbf{x} = \mathbf{y}$ ]	[ $\mathbf{x} = \mathbf{y}$ ]	[ $p$ ]
	{ }	{ }	{ $p$ }

Figure 4: Validation of Dead Store &amp; Alloca Elimination.

## Transformation

<b>B2</b>
6: $a = x + 1 \mapsto a = x + 1$

<b>B3</b>
0: $c = \phi(a@B2, b@B3) \mapsto t = \phi(x@B2, c@B3)$
1: <b>logical nop</b> $\mapsto c = t + 1$
2: $b = c + 1 \mapsto b = c + 1$

## Validation

<b>B2</b>
$\{ \}$ $\{ \}$ $\{ t, \hat{c} \}$
6: $a = x + 1 \mapsto a = x + 1$
$[ a = x + 1 ]$ $[ a = x + 1 ]$ $[ t, \hat{c} ]$
$\{ a = x + 1 \}$ $\{ \}$ $\{ t, \hat{c} \}$

<b>B3←B2</b>			
$\{ a = x + 1 \}$	$\{ \}$	$\{ t, \hat{c} \}$	
0: $c = \phi(a@B2, b@B3) \mapsto t = \phi(x@B2, c@B3)$			
$[ a = x + 1, c = a ]$	$[ t = x ]$	$[ t, \hat{c}, c ]$	
$\{ a = x + 1, c = a \}$	$\{ t = x \}$	$\{ t, \hat{c}, c \}$	
1: <b>logical nop</b> $\mapsto c = t + 1$			
$[ a = x + 1, c = a ]$	$[ t = x, c = t + 1 ]$	$[ t, \hat{c}, c ]$	
<b>trans_via_var</b> $c$ $a$ $(x + 1)$ <b>in</b> <b>src</b>			
$[ a = x + 1, c = a, c = x + 1 ]$	$[ t = x, c = t + 1 ]$	$[ t, \hat{c}, c ]$	
<b>replace_var</b> $t$ $x$ $c$ $(t + 1)$ $(x + 1)$ <b>in</b> <b>tgt</b>			
$[ a = x + 1, c = a, c = x + 1 ]$	$[ t = x, c = t + 1, c = x + 1 ]$	$[ t, \hat{c}, c ]$	
<b>remove_maydiff</b> $c$ $(x + 1)$			
$[ a = x + 1, c = a, c = x + 1 ]$	$[ t = x, c = t + 1, c = x + 1 ]$	$[ t, \hat{c} ]$	
$\{ \}$	$\{ \}$	$\{ t, \hat{c} \}$	
2: $b = c + 1 \mapsto b = c + 1$			
$[ b = c + 1 ]$	$[ b = c + 1 ]$	$[ t, \hat{c} ]$	
<b>intro_eq</b> $c$ $u$			
$[ b = c + 1, c = \hat{u} ]$	$[ b = c + 1, c = \hat{u} ]$	$[ t, \hat{c} ]$	
$\{ b = c + 1, c = \hat{u} \}$	$\{ c = \hat{u} \}$	$\{ t, \hat{c} \}$	
	$\{ b = c + 1, c = \hat{u} \}$		
<b>B3←B3</b>			
$\{ b = c + 1, c = \hat{u} \}$	$\{ c = \hat{u} \}$	$\{ t, \hat{c} \}$	
0: $c = \phi(a@B2, b@B3) \mapsto t = \phi(x@B2, c@B3)$			
$[ b = \hat{c} + 1, \hat{c} = \hat{u}, c = b ]$	$[ c = \hat{u}, t = c ]$	$[ t, \hat{c}, c ]$	
$\{ b = \hat{c} + 1, \hat{c} = \hat{u}, c = b \}$	$\{ c = \hat{u}, t = c \}$	$\{ t, \hat{c}, c \}$	
1: <b>logical nop</b> $\mapsto c = t + 1$			
$[ b = \hat{c} + 1, \hat{c} = \hat{u}, c = b ]$	$[ \hat{c} = \hat{u}, t = \hat{c}, c = t + 1 ]$	$[ t, \hat{c}, c ]$	
<b>trans_via_var</b> $c$ $b$ $(\hat{c} + 1)$ <b>in</b> <b>src</b>			
$[ b = \hat{c} + 1, \hat{c} = \hat{u}, c = b, c = \hat{c} + 1 ]$	$[ \hat{c} = \hat{u}, t = \hat{c}, c = t + 1 ]$	$[ t, \hat{c}, c ]$	
<b>replace_var</b> $\hat{c}$ $\hat{u}$ $c$ $(\hat{c} + 1)$ $(\hat{u} + 1)$ <b>in</b> <b>src</b>			
$[ b = \hat{c} + 1, \hat{c} = \hat{u}, c = b, c = \hat{c} + 1, c = \hat{u} + 1 ]$	$[ \hat{c} = \hat{u}, t = \hat{c}, c = t + 1 ]$	$[ t, \hat{c}, c ]$	
<b>trans_via_var</b> $t$ $\hat{c}$ $\hat{u}$ <b>in</b> <b>tgt</b>			
$[ b = \hat{c} + 1, \hat{c} = \hat{u}, c = b, c = \hat{c} + 1, c = \hat{u} + 1 ]$	$[ \hat{c} = \hat{u}, t = \hat{c}, c = t + 1, c = \hat{u} ]$	$[ t, \hat{c}, c ]$	
<b>replace_var</b> $t$ $\hat{u}$ $c$ $(t + 1)$ $(\hat{u} + 1)$ <b>in</b> <b>tgt</b>			
$[ b = \hat{c} + 1, \hat{c} = \hat{u}, c = b, c = \hat{c} + 1, c = \hat{u} + 1 ]$	$[ \hat{c} = \hat{u}, t = \hat{c}, c = t + 1, c = \hat{u}, c = \hat{u} + 1 ]$	$[ t, \hat{c}, c ]$	
<b>remove_maydiff</b> $c$ $(\hat{u} + 1)$			
$[ b = \hat{c} + 1, \hat{c} = \hat{u}, c = b, c = \hat{c} + 1, c = \hat{u} + 1 ]$	$[ \hat{c} = \hat{u}, t = \hat{c}, c = t + 1, c = \hat{u}, c = \hat{u} + 1 ]$	$[ t, \hat{c} ]$	
$\{ \}$	$\{ \}$	$\{ t, \hat{c} \}$	
2: $b = c + 1 \mapsto b = c + 1$			
$[ b = c + 1 ]$	$[ b = c + 1 ]$	$[ t, \hat{c} ]$	
<b>intro_eq</b> $c$ $u$			
$[ b = c + 1, c = \hat{u} ]$	$[ b = c + 1, c = \hat{u} ]$	$[ t, \hat{c} ]$	
$\{ b = c + 1, c = \hat{u} \}$	$\{ c = \hat{u} \}$	$\{ t, \hat{c} \}$	
	$\{ c = \hat{u} \}$		
	$\{ t, \hat{c} \}$		

Figure 5: Validation of Fold-Phi-Bin Optimization.

**vali**  $\mathcal{P}_{\text{src}} \mathcal{P}_{\text{tgt}} N R (P, Q, D) :=$

1. check  $\mathcal{P}_{\text{src}} = \mathcal{P}_{\text{tgt}}$  except for the code of each function.
2. check, for every function  $f$ ,  
**vali\_fun**  $\mathcal{P}_{\text{src}} \mathcal{P}_{\text{tgt}} f N(f) R(f) (P(f), Q(f), D(f)) = \text{true}$

---

**postcond**  $\mathcal{P}_{\text{src}} \mathcal{P}_{\text{tgt}} i_{\text{src}} i_{\text{tgt}} (P, Q, D) :=$

$(\hat{s}_{\text{src}}, \hat{s}_{\text{tgt}}, s_{\text{src}}, s_{\text{tgt}}, h_{\text{src}}, h_{\text{tgt}}) \in \llbracket (P, Q, D) \rrbracket (\alpha, \beta_{\text{src}}, \beta_{\text{tgt}})$   
 $\wedge (h_{\text{src}}, h_{\text{tgt}}) \in \mathcal{H}(\alpha, \beta_{\text{src}}, \beta_{\text{tgt}})$

1. check  $i_{\text{src}}$  and  $i_{\text{tgt}}$  are consistent.

**Construct**  $\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}}$  such that  $\alpha \sqsubseteq \alpha'$

$\wedge (\hat{s}_{\text{src}}, \hat{s}_{\text{tgt}}, s_{\text{src}}, s_{\text{tgt}}, h_{\text{src}}, h_{\text{tgt}}) \in \llbracket (P, Q, D) \rrbracket (\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$   
 $\wedge (h'_{\text{src}}, h'_{\text{tgt}}) \in \mathcal{H}(\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$

2. let  $(P', Q', D') = \text{rename\_to\_log\_var\_add\_md } i_{\text{src}} i_{\text{tgt}} (P, Q, D)$ .

**Construct**  $s'_{\text{src}}, s'_{\text{tgt}}$  such that

$(\hat{s}'_{\text{src}}, \hat{s}'_{\text{tgt}}, s'_{\text{src}}, s'_{\text{tgt}}, h_{\text{src}}, h_{\text{tgt}}) \in \llbracket (P', Q', D') \rrbracket (\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$   
 $\wedge (h'_{\text{src}}, h'_{\text{tgt}}) \in \mathcal{H}(\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$

3. let  $(P'', Q'', D'') = \text{del\_heap\_prop } i_{\text{src}} i_{\text{tgt}} (P', Q', D')$ .

$(\hat{s}'_{\text{src}}, \hat{s}'_{\text{tgt}}, s'_{\text{src}}, s'_{\text{tgt}}, h'_{\text{src}}, h'_{\text{tgt}}) \in \llbracket (P'', Q'', D'') \rrbracket (\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$   
 $\wedge (h'_{\text{src}}, h'_{\text{tgt}}) \in \mathcal{H}(\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$

4. let  $(P''', Q''', D''') = \text{del\_md\_add\_prop } i_{\text{src}} i_{\text{tgt}} (P'', Q'', D'')$ .

$(\hat{s}'_{\text{src}}, \hat{s}'_{\text{tgt}}, s'_{\text{src}}, s'_{\text{tgt}}, h'_{\text{src}}, h'_{\text{tgt}}) \in \llbracket (P''', Q''', D''') \rrbracket (\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$   
 $\wedge (h'_{\text{src}}, h'_{\text{tgt}}) \in \mathcal{H}(\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$

5. return  $(P''', Q''', D''')$

**vali\_fun**  $\mathcal{P}_{\text{src}} \mathcal{P}_{\text{tgt}} f N R (P, Q, D) :=$

1. check  $\text{CFG}(\mathcal{P}_{\text{src}}(f)) = \text{CFG}(\mathcal{P}_{\text{tgt}}(f))$ .
2. check  $P_{\text{pre}@B_{\text{init}}} = \emptyset \wedge Q_{\text{pre}@B_{\text{init}}} = \emptyset$ .
3. For every branch  $B_i \rightarrow B_j \in \text{CFG}(\mathcal{P}_{\text{src}}(f))$ ,  
check  $P_{\text{post}@B_i} \supseteq P_{\text{pre}@B_j \leftarrow B_i} \wedge Q_{\text{post}@B_i} \supseteq$   
 $Q_{\text{pre}@B_j \leftarrow B_i} \wedge D_{\text{post}@B_i} \subseteq D_{\text{pre}@B_j \leftarrow B_i}$ .
4. For each line  $n$  of each block  $B_i$ , check the following.

- (a) let  $(i_{\text{src}}, i_{\text{tgt}}) = \text{get\_instrs } \mathcal{P}_{\text{src}} \mathcal{P}_{\text{tgt}} N n$ .

$(\hat{s}_{\text{src}}, \hat{s}_{\text{tgt}}, s_{\text{src}}, s_{\text{tgt}}, h_{\text{src}}, h_{\text{tgt}}) \in \llbracket (P_{n-1}, Q_{n-1}, D_{n-1}) \rrbracket (\alpha, \beta_{\text{src}}, \beta_{\text{tgt}})$   
 $\wedge (h_{\text{src}}, h_{\text{tgt}}) \in \mathcal{H}(\alpha, \beta_{\text{src}}, \beta_{\text{tgt}})$

- (b) let  $(P', Q', D') =$

**postcond**  $\mathcal{P}_{\text{src}} \mathcal{P}_{\text{tgt}} i_{\text{src}} i_{\text{tgt}} (P_{n-1}, Q_{n-1}, D_{n-1})$ .

**Construct**  $\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}}, \hat{s}'_{\text{src}}, \hat{s}'_{\text{tgt}}$  such that  $\alpha \sqsubseteq \alpha'$

$\wedge (\hat{s}'_{\text{src}}, \hat{s}'_{\text{tgt}}, s'_{\text{src}}, s'_{\text{tgt}}, h'_{\text{src}}, h'_{\text{tgt}}) \in \llbracket (P', Q', D') \rrbracket (\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$   
 $\wedge (h'_{\text{src}}, h'_{\text{tgt}}) \in \mathcal{H}(\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$

- (c) let  $(P'', Q'', D'') =$   
**apply\\_inf\\_rules**  $R (P', Q', D')$ .

$(\hat{s}'_{\text{src}}, \hat{s}'_{\text{tgt}}, s'_{\text{src}}, s'_{\text{tgt}}, h'_{\text{src}}, h'_{\text{tgt}}) \in \llbracket (P'', Q'', D'') \rrbracket (\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$   
 $\wedge (h'_{\text{src}}, h'_{\text{tgt}}) \in \mathcal{H}(\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$

- (d) check  $P'' \supseteq P_n \wedge Q'' \supseteq Q_n \wedge D'' \subseteq D_n$ .

$(\hat{s}'_{\text{src}}, \hat{s}'_{\text{tgt}}, s'_{\text{src}}, s'_{\text{tgt}}, h'_{\text{src}}, h'_{\text{tgt}}) \in \llbracket (P_n, Q_n, D_n) \rrbracket (\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$   
 $\wedge (h'_{\text{src}}, h'_{\text{tgt}}) \in \mathcal{H}(\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$

Figure 6: Algorithm and Specification of SIMPLEBERRY.

**Unary Invariant** A unary invariant is a property  $P$  for the source or  $Q$  for the target program. What we define here applies to both the source and the target programs. The type of  $\llbracket P \rrbracket$  is as follows:

$$\llbracket P \rrbracket (\beta) \subseteq \text{Store} \times \text{Store} \times \text{Heap}$$

The semantics of  $P$  relies on isolated addresses  $\beta$  and relates the logical store  $\hat{s}$ , the physical store  $s$ , and the heap  $h$ , where the logical store is not connected to the physical state, and just maps logical variables to values.

The semantics of  $P$  is the conjunction of that of each element in  $P$ :

$$(\hat{s}, s, h) \in \llbracket P \rrbracket (\beta) \iff \forall p \in P. (\hat{s}, s, h) \in \llbracket p \rrbracket (\beta) .$$

The semantics of each basic property is given as follows:

- $\llbracket x = e \rrbracket (\beta) = \{(\hat{s}, s, h) \mid \llbracket x \rrbracket_{\hat{s}, s} = \llbracket e \rrbracket_{\hat{s}, s}\}$
- $\llbracket *p = x \rrbracket (\beta) = \{(\hat{s}, s, h) \mid \llbracket x \rrbracket_{\hat{s}, s} = h(\llbracket p \rrbracket_{\hat{s}, s})\}$
- $\llbracket p \neq q \rrbracket (\beta) = \{(\hat{s}, s, h) \mid \llbracket p \rrbracket_{\hat{s}, s} \neq \llbracket q \rrbracket_{\hat{s}, s}\}$
- $\llbracket \text{alc}(p) \rrbracket (\beta) = \{(\hat{s}, s, h) \mid \exists v. h(\llbracket p \rrbracket_{\hat{s}, s}) = v\}$
- $\llbracket \text{isol}(p) \rrbracket (\beta) = \{(\hat{s}, s, h) \mid \llbracket p \rrbracket_{\hat{s}, s} \in \beta\}$ ,

where the variable  $x$  can be either a physical or a logical variable.  $\llbracket x \rrbracket_{\hat{s}, s}$  looks up the logical store  $\hat{s}$  or the physical store  $s$  depending on whether  $x$  is logical or physical. The semantics  $\llbracket e \rrbracket_{\hat{s}, s}$  of expression is structurally defined as usual. The rest of the semantics is straightforward.

**May-Diff Set** The semantics of a may-diff set  $D$  relies on the memory injection  $\alpha$  and relates the source and target stores.

$$\llbracket D \rrbracket (\alpha) := \{(\hat{s}_{\text{src}}, \hat{s}_{\text{tgt}}, s_{\text{src}}, s_{\text{tgt}}) \mid \forall x \notin D. \llbracket x \rrbracket_{\hat{s}_{\text{src}}, s_{\text{src}}} \stackrel{\alpha}{\simeq} \llbracket x \rrbracket_{\hat{s}_{\text{tgt}}, s_{\text{tgt}}}\} .$$

It simply says that every variable not in  $D$  should be equivalent modulo  $\alpha$ .

We say that the logical and physical stores and heaps in the source and the target satisfy an invariant  $(P, Q, D)$ , denoted:

$$(\hat{s}_{\text{src}}, \hat{s}_{\text{tgt}}, s_{\text{src}}, s_{\text{tgt}}, h_{\text{src}}, h_{\text{tgt}}) \in \llbracket (P, Q, D) \rrbracket (\alpha, \beta_{\text{src}}, \beta_{\text{tgt}}) ,$$

if they satisfy all the semantics of  $P$ ,  $Q$ , and  $D$ .

## 5 Algorithm and Verification of Simpleberry

### 5.1 Algorithm

In this section, we explain the high-level algorithm of SIMPLEBERRY, given in Figure 6. Ignore the specifications of the algorithm written in blue color for now.

Given the source and the target programs  $\mathcal{P}_{\text{src}}, \mathcal{P}_{\text{tgt}}$  and hints  $N$  for logical nops,  $R$  for inference rules, and  $(P, Q, D)$  for invariants, the validator **vali** checks two things. It first checks that the optimization only changed the code of each function, and then checks that the translation made in each function is valid.

The function validator **vali\_fun** takes a function name  $f$  as argument and locally checks the validity of the translation in  $f$ . It first checks that the control flow graphs of the function



$f$  in the source and the target are equal. Second, it checks that the pre condition of the initial block is maximally permissive (*i.e.*,  $P, Q = \emptyset$  and  $D$  can be arbitrary). Third, it checks that whenever there is a branch from  $B_i$  to  $B_j$ , the post condition of  $B_i$  implies the corresponding pre condition of  $B_j$ . Finally, it validates each line of the function code.

At each line  $n$  of block  $B_i$ , **SIMPLEBERRY** first obtains the current instructions  $i_{\text{src}}$  and  $i_{\text{tgt}}$  in the source and the target programs. Here, it uses the logical nop hint  $N$  and the obtained instruction may be the logical nop. Second, it calculates a post condition of the current instructions. Third, it applies the inference rules  $R$  to the post condition. Finally, it checks that the derived invariant implies the invariant at the next line.

Finally, we explain the post condition calculator **postcond** in more detail. First of all, it checks that the current instructions are consistent; otherwise, it fails at validation. Examples of inconsistency include violation of the implicit heap invariant, or raise of different side effects such as error or function call. Second, it renames a variable  $x$  to the logical variable  $\hat{x}$  in the invariant if the variable  $x$  gets updated by the current instruction. Moreover, it adds both  $x$  and  $\hat{x}$  to the may-diff set.

Third, it removes from the invariant the heap properties that may not hold any more. For example, if an unknown part of the heap is updated or a non-read only function is called, we drop all heap equalities of the form  $*p = x$  from the invariant. Finally, we remove from the may-diff set the variable updated with equivalent values in the source and the target, and add new properties that hold for the new state after executing the current instruction.

## 5.2 Verification of Core Modules

In this section, we discuss how we verify the core modules of **SIMPLEBERRY**: **vali\_fun** and **postcond**.

First of all, since the source and the target programs have the same control flow graph by item 1 of **vali\_fun**, we can validate the translation block by block. Furthermore, since we align the instructions of the source and the target programs by the nop hint  $N$ , we can validate the translation line by line.

What we need to show for each function  $f$  is three fold. First, we need to show that whenever a function  $f$  is invoked in the source and the target programs, the program states satisfy the pre condition of the function (*i.e.*,  $P_{\text{pre}@B_{\text{init}}}, Q_{\text{pre}@B_{\text{init}}}, D_{\text{pre}@B_{\text{init}}}$ ). This follows from item 2 of **vali\_fun** because  $P_{\text{pre}@B_{\text{init}}} = \emptyset$ ,  $Q_{\text{pre}@B_{\text{init}}} = \emptyset$  and the local variables of each function can be assumed to initially have equivalent values in the source and the target. Second, whenever jumping to a new block, the program states satisfy the pre condition of the block. This follows from item 3 of **vali\_fun** trivially. Finally, we show that at each line of a block, if the program states satisfy the pre condition of the line, then after execution of the instructions (or function calls), the updated program states satisfy the post condition of the line. This last step is the main part of **SIMPLEBERRY** and we give the detailed specifications of the algorithm in blue color in Figure 6 .

We walk through the algorithm and explain the high-level ideas of the verification. First at 4-(b), we assume that the current store and heaps  $s_{\text{src}}, h_{\text{src}}, s_{\text{tgt}}, h_{\text{tgt}}$  of the source and the target programs satisfy the implicit heap invariant  $\mathcal{H}$  and the pre condition  $(P_{n-1}, Q_{n-1}, D_{n-1})$  for some partitioning  $\alpha, \beta_{\text{src}}, \beta_{\text{tgt}}$  and logical stores  $\hat{s}_{\text{src}}, \hat{s}_{\text{tgt}}$ . Then the function **postcond** calculates a post condition  $(P', Q', D')$  that is satisfied by the program states  $s'_{\text{src}}, h'_{\text{src}}, s'_{\text{tgt}}, h'_{\text{tgt}}$  after execution. Also, if **postcond** succeeds, the heap invariant  $\mathcal{H}$  is guaranteed to hold. At 4-(c), by the correctness of inference rules, we have that the derived invariant  $(P'', Q'', D'')$  holds. By 4-(d), we can easily see that the post condition  $(P_n, Q_n, D_n)$  also holds.

The verification of **postcond** is the most interesting. Before item 1 of **postcond**, we assume that the current state satisfies  $(P, Q, D)$  and the heap invariant  $\mathcal{H}$  for some partitioning and logical stores. We now walk through each step of **postcond**.

1. We check that the instructions executed or the functions called preserve the heap invariant

$\mathcal{H}$ . For example, when  $i_{\text{src}} = (\mathbf{p} = \mathbf{alloca}(\mathbf{x}))$  and  $i_{\text{tgt}} = (\mathbf{q} = \mathbf{alloca}(\mathbf{y}))$ , SIMPLEBERRY checks that the size of allocation is the same by checking for some variable  $z \notin D$  that the source has equality  $x = z$  and the target has  $y = z$ . Then, we can construct

$$\alpha' = \alpha \uplus (b_{\text{src}}, b_{\text{tgt}}), \beta'_{\text{src}} = \beta_{\text{src}}, \beta'_{\text{tgt}} = \beta_{\text{tgt}}$$

for  $b_{\text{src}}, b_{\text{tgt}}$  the newly allocated blocks in the source and the target. It is easy to see that the new states satisfy the heap invariant  $\mathcal{H}(\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}})$ . In the function call case, as long as we show that the same functions are called with the same arguments in the source and the target, we can freely assume that the returned states satisfy the heap invariant  $\mathcal{H}$  for some extended  $\alpha'$ .

2. Since we rename a newly defined variable  $x$  to the logical variable  $\hat{x}$ , we have to move the old value of  $x$  to the logical store. By renaming the equivalence between the source and the target may not hold any more, we add the renamed variables to the may-diff set. For example, if the pre condition has equality  $x = 5$  and  $i_{\text{src}} = (\mathbf{x} = \mathbf{1})$ , then we construct  $\hat{s}'_{\text{src}} = \hat{s}_{\text{src}}[x \leftarrow s_{\text{src}}(x)]$ , so that the renamed equality  $\hat{x} = 5$  holds. Note that since  $P', Q'$  still contains the heap equations for the old heaps  $h_{\text{src}}$  and  $h_{\text{tgt}}$ , we still use the old heaps in the specification.
3. Since `del_heap_prop` removes all heap equations that may not hold any more in the updated heaps, we can show that the new heaps  $h'_{\text{src}}$  and  $h'_{\text{tgt}}$  satisfy  $(P'', Q'', D'')$ . For example, if we call non-readonly functions, they may change the heap completely, so that SIMPLEBERRY drops all heap equations from the invariant.
4. Finally, we show that the newly added equations hold for the new state and the variables dropped from the may-diff set are equivalent in the new states.

### 5.3 Simulation

We show the final goal of program refinement between the source and the target programs (*i.e.*, the behavior of the target program is included in that of the source) using the simulation relation technique. To this end, we first define a relation  $\sim$  for each function that relates the source and the target programs, and show that the relation  $\sim$  is a simulation relation w.r.t. the instrumented operational semantics (*i.e.*, taking into account the logical nops). From this, by a standard argument, it follows that the target refines the source w.r.t. the instrumented semantics, which in turn implies that the same holds w.r.t. the real semantics by a simple coinduction.

**Local Simulation Relation** Suppose that the validation succeeds (*i.e.*,  $\mathbf{vali} \mathcal{P}_{\text{src}} \mathcal{P}_{\text{tgt}} \ N \ R(P, Q, D) = \mathbf{true}$ ). Let  $\mathcal{C}_{\text{src}}, \mathcal{C}_{\text{tgt}}$  be the static configurations of  $\mathcal{P}_{\text{src}}$  and  $\mathcal{P}_{\text{tgt}}$  after loading. Let  $f$  be any function of the program,  $ps_{\text{src}}, ps_{\text{tgt}}$  be any parent stack frames. Then we define a relation  $\sim$  as follows:

$$\begin{aligned} \mathcal{S}_{\text{src}} \stackrel{\alpha}{\sim} \mathcal{S}_{\text{tgt}} &\iff \exists \hat{s}_{\text{src}}, \hat{s}_{\text{tgt}}, \beta_{\text{src}}, \beta_{\text{tgt}}. \\ &(h_{\text{src}}, h_{\text{tgt}}) \in \mathcal{H}(\alpha, \beta_{\text{src}}, \beta_{\text{tgt}}) \wedge \\ &(\hat{s}_{\text{src}}, \hat{s}_{\text{tgt}}, s_{\text{src}}, s_{\text{tgt}}, h_{\text{src}}, h_{\text{tgt}}) \in \llbracket (P_{pc}, Q_{pc}, D_{pc}) \rrbracket (\alpha, \beta_{\text{src}}, \beta_{\text{tgt}}). \end{aligned}$$

where  $\mathcal{S}_{\text{src}} = ((pc, f, s_{\text{src}}) :: ps_{\text{src}}, h_{\text{src}})$  and  $\mathcal{S}_{\text{tgt}} = ((pc, f, s_{\text{tgt}}) :: ps_{\text{tgt}}, h_{\text{tgt}})$ . Here we simply relates any program states with the same program counter of the same function that satisfy the invariant at the current program counter  $pc$ . Of course, the initial states of the source and the target programs after loading are related by  $\sim$  for some  $\alpha$  because the pre condition of the initial block is empty.

Then using the verification results of **vali** in the previous section, we can show that the relation  $\sim$  is a simulation relation in the following sense.

**Theorem 1** (Simulation). For any  $\alpha$  and  $\mathcal{S}_{\text{src}} \overset{\alpha}{\sim} \mathcal{S}_{\text{tgt}}$ , we have:

$$\forall \mathcal{S}'_{\text{src}} \forall \varepsilon. (\mathcal{C}_{\text{src}}, N_{\text{src}}) \vdash \mathcal{S}_{\text{src}} \rightarrow^\varepsilon \mathcal{S}'_{\text{src}} \implies \exists \mathcal{S}'_{\text{tgt}}. (\mathcal{C}_{\text{tgt}}, N_{\text{tgt}}) \vdash \mathcal{S}_{\text{tgt}} \rightarrow^\varepsilon \mathcal{S}'_{\text{tgt}} \wedge$$

1.  $(\forall \mathcal{S}'_{\text{tgt}} \forall \varepsilon. (\mathcal{C}_{\text{tgt}}, N_{\text{tgt}}) \vdash \mathcal{S}_{\text{tgt}} \rightarrow^\varepsilon \mathcal{S}'_{\text{tgt}} \implies \exists \mathcal{S}'_{\text{src}} \exists \alpha' \supseteq \alpha. (\mathcal{C}_{\text{src}}, N_{\text{src}}) \vdash \mathcal{S}_{\text{src}} \rightarrow^\varepsilon \mathcal{S}'_{\text{src}} \wedge (\mathcal{S}'_{\text{src}} \overset{\alpha'}{\sim} \mathcal{S}'_{\text{tgt}})) \vee$
2.  $(i_{\text{src}} = \mathbf{ret} \ v_{\text{src}} \wedge i_{\text{tgt}} = \mathbf{ret} \ v_{\text{tgt}} \wedge \llbracket v_{\text{src}} \rrbracket_{s_{\text{src}}} \overset{\alpha}{\simeq} \llbracket v_{\text{tgt}} \rrbracket_{s_{\text{tgt}}}) \vee$
3.  $(i_{\text{src}} = (x_{\text{src}} = \mathbf{call} \ p_{\text{src}} \ \vec{a}_{\text{src}}) \wedge i_{\text{tgt}} = (x_{\text{tgt}} = \mathbf{call} \ p_{\text{tgt}} \ \vec{a}_{\text{tgt}}) \wedge \llbracket p_{\text{src}} \rrbracket_{s_{\text{src}}} \overset{\alpha}{\simeq} \llbracket p_{\text{tgt}} \rrbracket_{s_{\text{tgt}}} \wedge \llbracket \vec{a}_{\text{src}} \rrbracket_{s_{\text{src}}} \overset{\alpha}{\simeq} \llbracket \vec{a}_{\text{tgt}} \rrbracket_{s_{\text{tgt}}} \wedge (\forall \alpha' \supseteq \alpha \forall \beta'_{\text{src}} \forall \beta'_{\text{tgt}} \forall h'_{\text{src}} \forall h'_{\text{tgt}} \forall v_{\text{src}} \forall v_{\text{tgt}}. (h'_{\text{src}}, h'_{\text{tgt}}) \in \mathcal{H}(\alpha', \beta'_{\text{src}}, \beta'_{\text{tgt}}) \wedge (v_{\text{src}} \overset{\alpha'}{\simeq} v_{\text{tgt}}) \wedge \implies (pc'_{\text{src}}, s_{\text{src}}[x_{\text{src}} \mapsto v_{\text{src}}], h'_{\text{src}}) \overset{\alpha'}{\sim} (pc'_{\text{tgt}}, s_{\text{tgt}}[x_{\text{tgt}} \mapsto v_{\text{tgt}}], h'_{\text{src}}))) .$

where  $i_{\text{src}}$  and  $i_{\text{tgt}}$  are the instructions at the current pc of states  $\mathcal{S}_{\text{src}}$  and  $\mathcal{S}_{\text{tgt}}$  and  $pc'$  is the next line of pc.

This theorem says that if the source program does not get stuck, the target program should proceed and satisfies one of the three conditions: (i) for any target step, there is a corresponding source step, which results in the states again related by  $\sim$  for some  $\alpha'$  extending  $\alpha$ ; (ii) the programs are about to return the same values; or (iii) the programs are about to call the same function with the same arguments and, for any returned heaps and returned values related by  $\alpha'$  extending  $\alpha$ , the returned states are related by  $\sim$  for  $\alpha'$ .

We strongly believe that we can show that the program refinement result is derivable from this theorem. We have verified this simulation result for SIMPLEBERRY in Coq and are currently proving the refinement result. See Section 7 for more details.

## 6 Evaluation

In this section, we evaluate our validator and confirm its extensibility and practicality. During the experiment, we used the machine with Linux 3.5 operating system equipped with a quad-core of Intel 2.8GHz box with 24GB of main memory. Since our validator is based on Vellvm with LLVM 3.0, it is also written in Coq 8.3pl1 and extracted code is compiled by OCaml 3.12.0.

File	Description (and Validated Optimizations)	#vali/#defd
InstructionCombining	Main of instruction combine pass + Common optimization functions (dead_elim, add_assoc)	2 / 32
InstCombineAddSub	Optimizations for add and sub instructions (add_signbit, add_zext_bool, add_onebit, add_shift, add_sub, add_sub2, add_mul_fold, add_const_not, add_select_zero, add_select_zero, sub_add, sub_onebit, sub_mone, sub_const_not, sub_zext_bool, sub_const_add, sub_remove, sub_remove2, sub_sdiv, sub_shl, sub_mul, sub_mul2)	22 / 46
InstCombineMulDivRem	Optimizations for mul, div, and rem instructions (mul_mone, mul_neg, mul_bool, mul_shl, div_sub_srem, div_sub_urem, div_zext, div_mone, rem_zext, rem_neg, rem_neg2)	11 / 39
InstCombineAndOrXor	Optimizations for bitwise binary instructions (and_zero, and_same, and_mone, or_xor)	4 / 55
InstCombineSelect	Optimizations for select instructions (select_trunc, select_add, select_const_gt)	3 / 13
InstCombineCompares	Optimizations for compare instructions (cmp_onebit, cmp_eq, cmp_ult)	3 / 31
InstCombineLoadStoreAlloca	Optimizations for memory access instructions (load_load, store_load)	2 / 24
InstCombinePHI	Optimizations for phi instructions (fold_phi_bin)	1 / 13
InstCombineCasts	Optimizations for cast instructions (trunc_onebit)	1 / 35
InstCombineShifts	Optimizations for shift instructions (shift_undef)	1 / 14
InstCombineCalls	Optimizations for call and invoke instructions	0 / 50
InstCombineSimplifyDemanded	Logic for simplifying instructions based on information for how they are used	0 / 42
InstCombineVectorOps	Optimizations for vector instructions	0 / 7
TOTAL		50 / 401

Table 1: Micro optimizations that we have validated.

Program	LOC <sup>4</sup>	Files	Comp	Gen	Vali	Hints	ParseErr
bc-1.06	13K	19	2.1s	2.3s	23.7s	860	0
tar-1.27	17K	67	5.8s	2.8s	31.4s	948	0
less-451	24K	35	3.9s	2.4s	22.9s	1119	0
make-4.0	28K	27	4.0s	7.5s	122.0s	2393	0
a2ps-4.14	52K	83	16.3s	11.4s	299.0s	2926	0
python-3.4.1	304K	221	79.7s	633.3s	12779.6s	33946	15

Table 2: Validation Time

---

<sup>4</sup>We counted only the number of lines of .c files without headers.

SIMPLEBERRY was extended to support 50 micro optimizations of the instruction combine pass at small cost. At the first stage of the development, we selected three representative micro optimizations, `add_assoc`, `load_load`, and `fold_phi_bin`, and built the infrastructure for SIMPLEBERRY to be able to validate the three optimizations. After building the infrastructure, we started to add more inference rules and update the hint converter in order to validate more optimizations. The extension from 3 to 50 optimizations took 3 person-weeks for validation of the 47 micro optimizations without verifying the inference rules, and 2 person-weeks for the correctness proof of 32 inference rules out of 63 rules. Among the unproven 31 inference rules, we simply did not prove 18 rules just because they involve tedious proofs about arithmetic operations, while they are intuitively obvious. The remaining 13 rules are unprovable under the current semantics of Vellvm due to an issue with undefined values. See Section 7 for more details about the issue.

As of now, SIMPLEBERRY covers 50 micro optimizations and they are categorized in Table 1. We strongly believe SIMPLEBERRY can be used to validate almost all micro optimizations of the instruction combine pass. Three of the categories have not been implemented for the following reasons.

- `InstCombineCalls`: it simplifies function calls to intrinsic functions by using preliminary knowledge of their semantics. However, the semantics of the intrinsic functions such as `pow` have not been axiomatized in Vellvm, so we left them as a future work.
- `InstCombineSimplyDemanded`: it simplifies instructions based on analysis results of demanded bits. In this case, we can validate the analysis results in SIMPLEBERRY by adopting inference rules about masking. However, we have not done it yet, because generating proper hints requires clear understanding of the analysis, which is not straightforward.
- `InstCombineVectorOpts`: it simplifies vector operations, but the vector type is not supported by Vellvm.

To show SIMPLEBERRY is practical, we applied it to our benchmarks and the results are presented in Table 2. `Comp`, `Gen`, `Vali`, and `Hints` denotes compile time, core hint generation time, validation time, and the number of generated hints, respectively. `ParseErr` denotes the number of source files that include unsupported syntax such as inline assembly and indirect branch. In our experiment, optimizations on 15 files were not validated for that reason. Except for the optimizations from the 15 files, all other 42000 generated hints were successfully validated in the experiment. The validation (`Gen+Vali`) took 5 to 168 times longer than the original compilation (`Comp`), but it can be parallelized as shown in Figure 1. Indeed, we performed the validations in parallel using 3 cores and added up the time consumed in the table above.

## 7 Discussion and Future Work

We discuss issues raised during our development of SIMPLEBERRY and future work.

**Issues with Formal Semantics of LLVM** In the course of the verification, we found several issues with both the formal semantics of Vellvm and the official semantics of LLVM. We discuss the issues and our opinion about them. Note that Vellvm provides two versions of formal semantics: deterministic one and non-deterministic one. The former is closer to the CompCert memory model and the latter is closer to the official semantics of LLVM. For our verification, we used the deterministic semantics because it is more natural than the other.

The deterministic semantics of Vellvm models the undefined value as a special value, called `undef`. This value is propagated through arithmetic operations. For example, we have `undef + 1 = undef`. The problem here is that we lose some information when such propagation occurs.

Due to this problem, we cannot prove the correctness of several optimizations in the instruction combine pass. As an example, consider the following optimization.

B1
5: $x = a - 1 \quad \mapsto \quad x = a - 1$
6: $y = a - x \quad \mapsto \quad y = a - x$
7: $z = y + 1 \quad \mapsto \quad z = 1 + 1$

Intuitively, since  $y = a - x = a - (a - 1) = 1$ , we can replace  $y$  by 1. However, if  $a$  is **undef**, then  $z = \mathbf{undef}$  in the source and  $z = 2$  in the target. This is problematic because, in the deterministic semantics, **undef** is interpreted as **false** in the branch while 2 is interpreted as **true**.

The official semantics of LLVM treats the undefined value as a set of all possible values. This makes the non-deterministic choice of undefined value as a lazy effect. For example, after executing  $x = \mathbf{undef}$ , if you read  $x$  twice, you may get different values. The purpose of this semantics is to allow to freely replace  $x$  by **undef** when you have  $x = \mathbf{undef}$ . However, this semantics leads to a somewhat weird optimization. Here is an example.

```

1: unsigned char bar(unsigned char *i) {return *i;}
2: unsigned char foo() {
3:   unsigned char i;
4:   bar(&i);
5:   return i;
6: }
7: int main() {
8:   unsigned int x = (unsigned int) foo();
9:   printf(x < 2 ? "a" : "b");
10:  printf(x > 1 ? "a" : "b");
11:  return 0;
12: }
```

If you compile this program using clang with optimization level 2 (*i.e.*, -O2), the compiled program prints out **bb**. We tested this program with clang version 3.0 and 3.5.0. However, this optimization makes debugging hard because if you look at the code at line 9 and 10, you get contradiction because  $x \geq 2$  and  $x \leq 1$  should hold. Though not very clear, the standard C semantics seems to allow this because the program can go wrong at line 9 due to the uninitialized of the variable  $i$  at line 3. The reason why “bb” is printed is because LLVM propagates the effect of undefinedness raised for the variable  $i$  at line 3 to the variable  $x$  at lines 9 and 10 (*i.e.*, replacing  $x$  by **undef**). Then both **undef** < 1 and **undef** > 0 are optimized to 0 to take the else branch.

Our plan is to model **undef** as an immediate effect of non-deterministic choice (not a special value nor a lazy effect). This simplifies the formal semantics and thus makes verification a lot easier. For example, we can have algebraic properties such as  $x = a - 1 \wedge y = a - x \implies y = 1$ , which is not the case in the presence of **undef** as discussed above. A down side is that we cannot allow propagating the **undef** command. We think this is not too bad because **undef** implicitly means that something went wrong and optimizations in an erroneous state might not be so useful.

There is another issue with LLVM optimization. The LLVM dead code eliminator removes a read-only function with even divergence effect. If you compile the example below using clang with optimization level 1 (*i.e.*, -O1), both in version 3.0 and 3.5.0, the compiled program prints out the “boom!” message.

```

1: void foo() { while (1) { } }
```

```

2: void bar() {
3:   foo();
4:   for (int i = 0; i < 100; ++i) { }
5: }
6: int main() {
7:   bar();
8:   printf("boom!\n");
9: }

```

Intuitively this code should not terminate without printing “boom!” because `foo()` diverges. The standard C semantics also does not allow this but for a little subtle reason. The C semantics basically allows to optimize away a benign infinite loop unless the loop condition is a constant. In the above example, since the loop condition at line 4 is a constant, it should not be removed.

Regardless of whether an infinite loop has a constant condition or not, optimizing away any infinite loop does not conform to the standard notion of program refinement and thus we cannot prove such a dead code optimization.

One possible direction is to provide a flag like `-safe` that prevents performing unsafe optimizations such as undef propagation and read-only function elimination unless we make sure that they are correct w.r.t. the formal semantics (*e.g.*, the read-only function always terminates). Another direction is to find a reasonable formal semantics that allows a benign infinite loop elimination.

Finally, LLVM has a notion of poison value, which is an undefined value with additional flag that it is poisonous. This value does not appear in the program syntax, but only in the dynamic semantics. We leave how to formalize poison values as a future work.

### Current Status of Our Development

We discuss the current status of development of SIMPLEBERRY.

First, we fully verified Theorem 1 in Coq 8.3 pl1 for SIMPLEBERRY assuming that the source and the target programs are well-formed and that functions with `readonly` tag do not modify the global and the parent’s heap. The theorem in Coq corresponding to Theorem 1 requires the well-formedness of the source and the target program, and moreover, when the optimization relies on a function’s `readonlyness`, the correctness of `readonly` tags.

Currently, we do not provide such properties and thus, technically speaking, our proof is incomplete. However, we can provide proofs for such properties by running a well-formedness checker and a `readonlyness` checker when necessary. Since well-formedness checking is syntax-direct and easy, we plan to develop such a checker in Coq and prove it correct. On the other hand, since `readonlyness` checking is quite complex, we plan to use the LLVM’s `readonlyness` analyzer and validate its result by verified validation. We expect such additional checking would not take so long because well-formedness checking is very fast and, on the other hand, rather slow `readonlyness` checking is not invoked very often.

We also have memory fault bugs in our hint converter due to the foreign function interface (FFI) with C. In our hint converter, we use the LLVM bitcode parser for OCaml included in the official release, which invokes a parser written in C via FFI. We experience memory fault bugs when using the parser such as segmentation fault, out-of-memory or even divergence. We temporarily removed all such buggy behavior by inserting explicit GC calls, or strangely by enabling the debug flag in the parser. This is not a proper solution and we plan to properly debug and fixed the bugs.

Finally, there are several features of LLVM that Vellvm do not formalize. Among those, we only discuss features that appear during compilation of C source files. First, the `switch` instruction is not formalized and thus we remove it by running the `lower-switch` pass. Second, the `inline assembly` and the `indirect branch` instructions are not formalized and thus we could not validate any source files using such features. Third, Vellvm do not have concrete

semantics for known functions like `pow` and thus we did not validate related optimizations such as optimizing `pow(x, 0)` to 1. We can easily support it by axiomatizing the behavior of such functions. Fourth, Vellvm ignores flags to instructions that can raise poison values, such as `nsw` and `inbounds` flags.

Formalizing all these features is an interesting future work.

## 8 Related Work

Compiler verification has been a grand challenge [2] and there has been a great amount of work that made such a goal feasible.

**Verified Compilation** The most important step towards compiler verification is Leroy’s seminal work on the CompCert compiler [4, 8]. It is a fully verified C compiler written in Coq. Despite of various efforts, there has been no bug found for verified parts of the compiler [11, 3]. Our LLVMBERRY project is strongly motivated by this result.

In particular, the CompCert memory model has a great influence on other compiler verification work. For example, Vellvm’s memory model is also based on the CompCert model. We also used several lemmas from the CompCert project for our verification of SIMPLEBERRY.

**LLVM Formalization and Verification** The Vellvm framework [12] provides a formalization of LLVM IR and our work is directly built up on this framework. Vellvm has been also used to verify the `mem2reg` pass of the LLVM compiler [13]. However, for this verification, the `mem2reg` pass is rewritten and verified in Coq. The reason for re-implementation is that the original `mem2reg` pass temporarily breaks the well-formedness of LLVM IR.

We believe, however, that it is possible to develop a verified validator for the original `mem2reg` pass because, in verified validation, we do not need to show the preservation of well-formedness. Instead, we just need to check the well-formedness using a well-formedness checker whenever the well-formedness is recovered. In other words, it is fine for well-formedness to be broken temporarily as long as we do not rely on the property at that point.

We believe this because indeed SIMPLEBERRY does not rely on the full SSA property, but only on the well-typedness.

**Verified Validation** Tristan and Leroy [8, 9, 10] first applied the verified validation approach to several compiler optimizations such as instruction scheduling, lazy code motion, and software pipelining. In this work, they have shown that verified validation can be much more effective than verified compilation for some optimizations.

Barthe et al. [1] presented a formally verified SSA-based, middle-end compiler on CompCert. In this work, they applied the verified validation to a complex global value numbering (GVN) optimization. This is another example where verified validation is more effective than verified compilation.

**Translation Validation** The translation validation approach is to develop a general-purpose algorithm that checks the validity of compilation results, which motivated the verified validation approach. Necula [6] developed an algorithm using symbolic evaluation and Tristan et al. [7] developed a translation validator based on a value-graph algorithm and applied it to several optimization passes in LLVM.

A downside of this approach is that the algorithm is not complete. For example, the validator developed by Tristan et al. [7] has not been applied to the instruction combine pass. More importantly, since the validator is not verified, we cannot trust the validation result with the highest confidence.



## References

- [1] G. Barthe, D. Demange, and D. Pichardie. A formally verified SSA-based middle-end: Static single assignment meets Compcert. In *ESOP*, 2012.
- [2] T. Hoare. The verifying compiler: A grand challenge for computing research. *J. ACM*, 50(1):63–69, Jan. 2003.
- [3] V. Le, M. Afshari, and Z. Su. Compiler validation via equivalence modulo inputs. In *PLDI*, 2014.
- [4] X. Leroy. Formal verification of a realistic compiler. *Communications of the ACM*, 52(7):107–115, 2009.
- [5] X. Leroy and S. Blazy. Formal verification of a C-like memory model and its uses for verifying program transformations. *J. Autom. Reason.*, 41(1):1–31, July 2008.
- [6] G. C. Necula. Translation validation for an optimizing compiler. In *PLDI*, 2000.
- [7] J.-B. Tristan, P. Govereau, and G. Morrisett. Evaluating value-graph translation validation for LLVM. In *PLDI*, 2011.
- [8] J.-B. Tristan and X. Leroy. Formal verification of translation validators: A case study on instruction scheduling optimizations. In *POPL*, 2008.
- [9] J.-B. Tristan and X. Leroy. Verified validation of lazy code motion. In *PLDI*, 2009.
- [10] J.-B. Tristan and X. Leroy. A simple, verified validator for software pipelining. In *POPL*, 2010.
- [11] X. Yang, Y. Chen, E. Eide, and J. Regehr. Finding and understanding bugs in C compilers. In *PLDI*, 2011.
- [12] J. Zhao, S. Nagarakatte, M. M. Martin, and S. Zdancewic. Formalizing the LLVM intermediate representation for verified program transformations. *SIGPLAN Not.*, 47(1):427–440, Jan. 2012.
- [13] J. Zhao, S. Nagarakatte, M. M. Martin, and S. Zdancewic. Formal verification of SSA-based optimizations for LLVM. In *PLDI*, 2013.