

Extensible Verified Translation Validation for LLVM

Abstract

Both efficiency and reliability are important for compilers. However, mainstream compilers tend to focus more on efficiency than on reliability. They are complex pieces of software usually written in C/C++ performing various optimizations, thereby containing many bugs. Indeed, recent random testing tools found hundreds of bugs in GCC and LLVM.

To increase the reliability of mainstream compilers without sacrificing their performance, we propose *extensible verified translation validation*. The key ideas are four-fold. First, we perform translation validation. Since we compile first and validate the result later, we do not sacrifice any compilation performance. Second, we modify the compiler itself to generate validation proofs of translations. This way we can always produce correct explicit proofs, or identify compiler bugs. Third, we design an extensible validator that validates the proofs given by the compiler. Our validator is based on relational Hoare logic and can be easily specialized to a particular optimization by simply adding new inference rules. Finally, we can also verify the validator and inference rules in order to reduce the trusted computing base.

We apply our approach to validating 146 micro-optimizations of the `instcombine` pass, the main GVN-PRE algorithm of the `gvn` pass, and register promotion in the `sroa` pass in LLVM 3.7.1. We found three compiler bugs (two fixed, one confirmed) during the development. We ran our validator for SPEC CINT2006 benchmarks and other five projects written in C, with 3.9M lines of code in total. We succeeded at validating 1M translation steps, except for 254 failures, all of which are due to the three compiler bugs we found. We also formally verified correctness of the core part of the validator in the Coq proof assistant.

1. Introduction

Mainstream compilers such as GCC and LLVM are complex pieces of software satisfying various requirements. They are required to generate very efficient code and thus perform many complex optimizations. Also the compilers themselves are required to be very efficient (*i.e.*, to consume less time and memory during compilation) and thus are usually written in C/C++ using sophisticated data structures.

Due to this complexity, it is hard to make mainstream compilers very reliable, which may be seen as one of the most important requirements. Indeed, like other large complex software, current mainstream compilers have many bugs. For

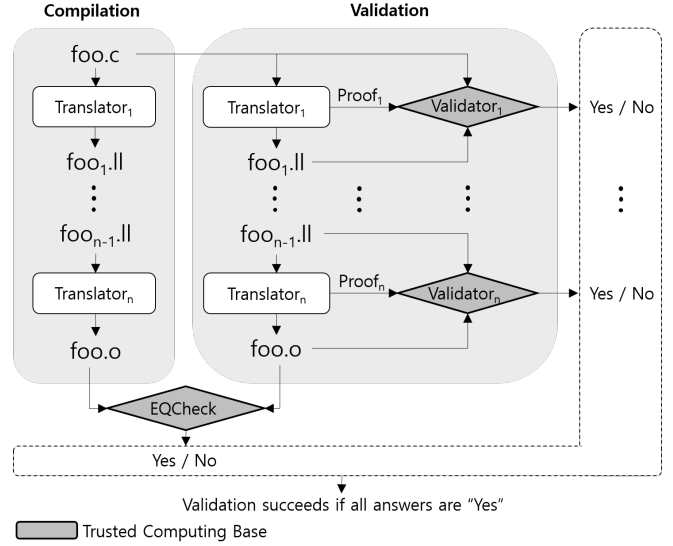


Figure 1. Framework

example, recent random testing tools such as CSmith [24, 8] and EMI [12] have found hundreds of bugs in GCC and LLVM.

In order to achieve a high level of reliability without sacrificing performance, we propose a practical approach to compilation verification, which we call *extensible verified translation validation*, and demonstrate its applicability to mainstream compilers. We apply our approach to the global value numbering, the partial redundancy elimination, the register promotion, and 146 micro-optimizations of the instruction combining optimizations in LLVM 3.7.1.

In order to achieve our goal, we design our framework by combining the best ideas from existing approaches (see §1.1) and develop an *extensible relational Hoare-logic (ERHL)* validator to implement the framework (see §1.2).

1.1 Our Framework

We briefly introduce our framework, which is a combination of (verified) translation validation [20, 19, 21, 23, 22, 18] and (foundational) proof carrying code [16, 6] (see §7.1 for comparison).

Problems The main problem with verifying mainstream compilers is simply that they are too large to directly verify. For example, GCC and LLVM consist of million lines of code. Our solution for this problem is to follow the translation

validation approach [20, 19]. The idea is that, instead of verifying the compiler itself, we validate the correctness of compilation results.

Translation Validation with Proofs In our framework, shown in Figure 1, we separate compilation and validation phases. For compilation, depicted in LHS of Figure 1, we use the original compiler to translate, for example, the source code `foo.c` to the target code `foo.o`, which is internally translated through intermediate representations `foo1.ll, ..., foon-1.ll`. At any time after the compilation, we can conduct post validation, depicted in RHS of Figure 1. For this, we first run the original compiler modified to additionally generate proofs of correctness, `Proofi`. Then the validators `Validatori` check whether `fooi-1.ll` is correctly translated to `fooi.ll` using `Proofi`.

Finally, we succeed at validation of the original compilation if each validation conducted by `Validatori` succeeds and the final result of the original compilation (`foo.o` in LHS) coincides with that of the compilation for validation (`foo.o` in RHS).

Verification of Validators In our framework, the trusted computing base (TCB) includes only the validators `Validatori` and the equality checker `EQCheck`. In particular, the proof generation code in the modified compiler is not a part of TCB because any incorrect proof would be invalidated by the validator.

Optionally, we can further reduce the TCB by verifying the correctness of validators: whenever validation of a translation succeeds, the translation is semantics-preserving. Indeed, we implement our validator in the Coq proof assistant and prove correctness of its core using the formal semantics of LLVM IR (Intermediate Representation) from the Vellvm project [26]. In this case, the TCB reduces to the validity of the formal semantics and the correctness of the Coq proof checker.¹

Advantages The advantages of our framework is four-fold. First, our framework does not sacrifice compilation performance because we completely separate compilation from validation. Second, although the post validation may take long, the validation time can be reduced dramatically by parallelizing the validation jobs because they are completely independent. Third, we can always produce correct explicit proofs of translation, unless there is a compiler bug, since we can get enough information from the compiler. Finally, we can verify the validators to achieve the highest level of reliability, as in fully verified compilers such as CompCert [13]

1.2 Extensible Relational Hoare-Logic Validator

Now the main question remaining is how to generate and validate proofs.

¹ In our experiment, we did not target the frontend translation from `foo.c` to `foo1.ll` and the backend translation from `foon-1.ll` to `foo.o`.

Problems A possible solution would be to follow the verified validation approach [21, 23, 22, 18]. For each translator `Translatori`, we develop a highly specialized validator `Validatori` that fully understands how `Translatori` works and validates the correctness of its translation results, when necessary, using some helper information `Proofi` given by `Translatori`.

A problem with this approach is that it is very costly to develop a new validator for each translation pass in mainstream compilers. Typically a mainstream compiler such as GCC and LLVM includes over a hundred major translation passes. Things get even worse if we want to verify the validators: we have to verify over a hundred validators!

Our Approach To solve this problem, we take a different approach. We develop a *single* general validator that can validate correctness proofs given by different translation passes. Here, the general validator is a simple proof checker whereas the compiler provides detailed proofs.

Our general validator, called *extensible relational Hoare-logic (ERHL)* validator, can be applied to translation passes that do not change the control flow graph. For those that change the control flow graph such as loop unrolling or loop unswitching, we need to develop a more general validator, which we leave as future work.

To give a high level idea, we briefly illustrate what constitutes a proof and how it is validated in our ERHL system. For this, suppose we validate the translation from `src.ll` to `tgt.ll` with `Proof`. Also, for simplicity of presentation, suppose the source and target programs consist of a sequence of instructions without any branch.

$$\begin{array}{ccc}
 [\text{src.ll}] & & [\text{tgt.ll}] \\
 \vdots & & \vdots \\
 & P_{i-1} & \\
 \text{inst}_i^{\text{src}} & & \text{inst}_i^{\text{tgt}} \\
 & Q_i & \\
 & P_i & \\
 \text{inst}_{i+1}^{\text{src}} & & \text{inst}_{i+1}^{\text{tgt}} \\
 \vdots & & \vdots
 \end{array}$$

1. **Instruction Alignment** Proof inserts logical no-op instructions to align the instructions of `src.ll` and `tgt.ll` one by one as shown above. This is possible because the translation does not change the control flow.
2. **Invariant Generation** Proof provides a relational invariant P_i between program states of `src.ll` and `tgt.ll` after each pair $(\text{inst}_i^{\text{src}}, \text{inst}_i^{\text{tgt}})$ of source and target instructions.
3. **Post-Invariant Computation** The validator (i) checks that P_0 is the default invariant that should hold for all initial states; (ii) checks that $\text{inst}_i^{\text{src}}$ and $\text{inst}_i^{\text{tgt}}$ produce the same observable events (e.g., calling the same function with the same arguments) under all pairs of program states satisfying P_{i-1} ; and (iii) computes a post-invariant

Q_i that should hold for all program states given after executing $inst_i^{src}$ and $inst_i^{tgt}$ under any program states satisfying P_{i-1} . This validator is completely independent from the translation passes and tries to generate as strong post-invariants as possible.

4. Inference Rule Application Finally, Proof shows the validator how to prove Q_i implies P_i (i.e., any pair of program states satisfying Q_i also satisfies P_i). For this, we pre-install, in the validator, *inference rules* possibly specific to particular translation passes. Then Proof provides the validator with a sequence of inference rules to apply in order to infer P_i from Q_i .

We can easily see that the source and target programs produce the same observable events if the validation succeeds.² Here the TCB includes the correctness of the validator (mainly, that of the post-invariant computation) and that of the inference rules added. We can optionally verify them to reduce the TCB.

Applicability Our validator is easily extensible: for supporting a new optimization, you just need to install new inference rules possibly specialized for the optimization and verify them. It costs much lower than implementing and verifying a specialized validator for the optimization.

We believe our approach applies to most of the optimizations that do not change control flow graph, which comprises 66.7% of the -O2 optimization passes of LLVM 3.7.1. This is because the majority of compiler optimizations in LLVM rely on relatively simple (flow-sensitive) analysis, which we can handle in a relational Hoare logic.

1.3 Result

We developed a prototype ERHL validator for the LLVM compiler, and successfully validated compilations of 146 micro-optimizations of the `instcombine` pass, the main GVN-PRE algorithm of the `gvn` pass, and register promotion in the `sroa` pass for SPEC CINT2006 C benchmarks [5] and other five C projects, with 3.9M lines of code in total. In the benchmark we successfully validated 1,191,063 translations, except for 254 failures that are due to the three compiler bugs we found during this project, and 197,652 (16.6%) that are currently not supported in the validator, for example, because they contain instructions that are not formalized yet. We also prove correctness of the core part of the validator in Coq, though leaving admits on minor details.

2. Overview

In this section, we present more details of how our validation system works using the add-*assoc* optimization as a motivating example, which is a part of the `instcombine` pass of the LLVM compiler.

$\begin{array}{c} : \\ 10: x := \text{add } a \ 1 \\ : \\ 20: y := \text{add } x \ 2 \\ : \end{array}$	\rightsquigarrow	$\begin{array}{c} : \\ 10: x := \text{add } a \ 1 \\ : \\ 20: y := \text{add } a \ 3 \\ : \end{array}$
--	--------------------	--

Figure 2. An instance of the add-*assoc* optimization

Algorithm 1 AddAssoc (P:program)

```

1: for  $y$  in  $\text{Reg}(P)$  do
2:   match  $\text{FindDef}(P, y)$  with
3:   |  $\text{Some}(l_2, \text{add}(\text{reg } x) (\text{const } C_2)) \Rightarrow$ 
4:     match  $\text{FindDef}(P, x)$  with
5:     |  $\text{Some}(l_1, \text{add}(\text{reg } a) (\text{const } C_1)) \Rightarrow$ 

Propagate( $\text{src}, x = \text{add } a \ C_1, l_1, l_2$ )
          Infrule( $\text{src}, \text{assoc}(\text{add}, x, y, a, C_1, C_2), l_2$ )
          Infrule( $\text{reduce\_maydiff}(y), l_2$ )


6:        $C = \text{Simplify}(\text{add}, C_1, C_2)$ 
7:       ReplaceAt( $P, l_2, y = \text{add}(\text{reg } a) (\text{const } C)$ )
8:     |  $\_ \Rightarrow ()$  end match
9:   |  $\_ \Rightarrow ()$  end match
10: end for

```

2.1 The add-*assoc* Optimization

Figure 2 shows an example of the add-*assoc* optimization. This gives the opportunity to remove the instruction `10: x := add a 1` later, if x is not used anymore. Correctness of this translation is quite clear intuitively. Thanks to the SSA (Static Single Assignment) property [9]³, there are no intervening writes to a and x between line 10 and 20, and thus at line 20, $x = a + 1$ holds, from which $y = x + 2 = (a + 1) + 2 = a + 3$ follows.

We give a simplified version of the compiler code for add-*assoc*, presented in a functional style for presentation purpose, in Algorithm 1. The code in the dashed box is what we add to generate proofs for validation, which we will discuss later. The algorithm proceeds as follows:

Lines 1-3 For a source program P , we find an instruction of the form $l_2: y := \text{add } x \ C_2$ with C_2 constant. Note that $\text{Reg}(P)$ is the set of all registers defined in P and $\text{FindDef}(P, y)$ finds the instruction defining the register y , which is unique thanks to the SSA property. For e.g., in Figure 2, we can find `20: y := add x 2`.

Lines 4-5 We check if x is defined by an instruction of the form $l_1: x := \text{add } a \ C_1$ with C_1 constant. For e.g., in Figure 2, the register x is defined by `10: x := add a 1`.

Lines 6-7 If that is the case, we compute the constant $C = \text{add } C_1 \ C_2$ and replace the instruction at l_2 with $y := \text{add } a \ C$. For e.g., in Figure 2, `20: y := add x 2` is replaced by `20: y := add a 3`.

² Technically, we show behavioral refinement instead of equivalence. For simplicity of presentation, we use behavioral equivalence for now.

³ The SSA property is that for every register, there is exactly one instruction defining it (i.e., assigning a value to it).

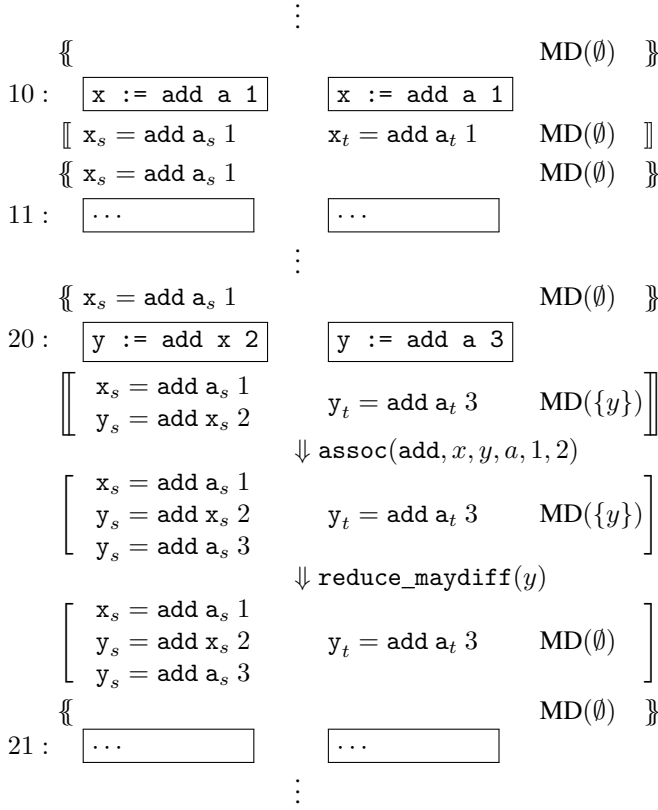


Figure 3. Validation of an add-assoc proof in ERHL

2.2 Generation and Validation of Translation Proof

Now we generate and validate proofs for add-assoc translations in ERHL, as outlined in the introduction. To help understanding, [Figure 3](#) shows the proof and its validation for the example translation of [Figure 2](#). First, we do not need to insert any logical no-ops because the source and target programs are already well aligned.

2.2.1 Invariant Generation

Invariant Format Before we proceed, we show the format of invariants in ERHL. An ERHL invariant is a triple (S, T, M) , where S is an invariant that should hold for the source state; T is an invariant for the target state; and M is an invariant relating the source and target states.

First, every ERHL invariant implicitly impose that the source and target memories should be equivalent. More precisely, we use the CompCert-style memory injection relation [14] (see [§3.5](#) for more details).

Second, the source and target invariants, S and T , are sets of propositions that should hold for the source and target states, each of which can contain various forms of predicates. Though we only use the equality predicate for add-assoc, we will see more predicates for other optimizations in [§3](#). As examples, $\{x_s = \text{add } a_s \ 1, y_s = \text{add } x_s \ 2\}$ is a source invariant and $\{y_t = \text{add } a_t \ 3\}$ is a target invariant. Here and

henceforth, x_s and x_t represent the value of the register x in the source and target states, respectively.

Finally, the invariant M is a set of registers, called *maydiff set*, that may contain different values in the source and target states. In other words, all the registers not in M should have the same value in the source and the target states, which we denote by $\text{MD}(M)$:

$$\text{MD}(M) \iff \forall x \notin M. x_s = x_t.$$

A main advantage of ERHL invariants is that we can use the standard algorithm of Hoare logic to compute post invariants because ERHL invariants are mainly unary (only maydiff sets are relational). Though mainly unary, ERHL invariants can encode more general forms of relational invariants (see [§3.3](#)), which we believe are sufficiently expressive to handle standard compiler optimizations (preserving the control flow).

Invariant Generation As shown in [Figure 3](#), we generate an appropriate ERHL invariant at each line (see those in $\{\{\}\}$). In ERHL, the default invariant is $(\emptyset, \emptyset, \emptyset)$, which denotes that the source and target memories are equivalent and the source and target registers contain the same values.

To prove correctness of add-assoc, we need to know $x_s = \text{add } a_s \ C_1$ at the point of defining y , and thus we add $x_s = \text{add } a_s \ C_1$ to the source invariants between line l_1 , where x is defined, and line l_2 , where y is defined. The code `Propagate(src, x = add a C1, l1, l2)` in the dashed box of [Algorithm 1](#) generates such additional invariants. More precisely, in order to reduce the proof size, the compiler just outputs the command `Propagate(src, x = add a C1, l1, l2)`, from which a preprocessor in the validator, called *invariant propagator*, generates the full invariants. Note that the invariant propagator is not a part of the TCB.

2.2.2 Post-Invariant Computation

As outlined in the introduction, the ERHL validator proceeds in three steps. First, it checks whether the initial invariant is the default invariant $(\emptyset, \emptyset, \emptyset)$, which is the case in [Figure 3](#). Second, it checks whether the same observable events are produced at each line. It is the case in [Figure 3](#) because (i) at line 20, no observable events are produced; and (ii) at the other lines, the source and target instructions are identical and the maydiff sets are empty implying that the source and target states are equivalent.

Finally, at each line in [Figure 3](#), the validator computes a post-invariant, that in $\{\{\}\}$ just after the line, from the pre-invariant, that in $\{\{\}\}$ just before the line. The post-invariant computation algorithm is simple. It basically (i) removes from the source and target invariants the propositions that may be invalidated by the instructions executed; (ii) adds to the source and target invariants new propositions induced from the instructions executed; and (iii) adds to the maydiff set the registers just defined if the instructions executed are not identical, or their arguments are in the maydiff set. In

Figure 3, for example, at line 10 we add $x_s = \text{add } a_s \ 1$ and $x_t = \text{add } a_t \ 1$; and at line 20 we add $y_s = \text{add } x_s \ 2$, $y_t = \text{add } a_t \ 3$ and $\text{MD}(\{y\})$.

2.2.3 Inference Rule Application

We finally infer, at each line, the invariant in $\{\}$ given by the compiler, from that in $[\]$ computed by the validator.

Except line 20, the inference is done automatically by the validator because it is trivial. In particular, since there are no intervening writes to a and x between line 10 and 20 thanks to the SSA property, the invariant $x_s = \text{add } a_s \ 1$ is never removed by the post-invariant computation algorithm, so that the invariant trivially holds at the following line.

At line 20, the inference is guided by the proof from the compiler. First, the code $\text{Infrule}(\text{src}, \text{assoc}(\text{add}, x, y, a, C_1, C_2), l_2)$ in the dashed box of Algorithm 1 outputs a command that applies the following pre-installed assoc rule with the given arguments at line l_2 .

$$\frac{(\text{assoc}(\text{add}, x, y, a, C_1, C_2)) \quad x = \text{add } a \ C_1 \quad y = \text{add } x \ C_2 \quad C = C_1 + C_2}{y = \text{add } a \ C}$$

In Figure 3, we have $y_s = \text{add } a_s \ 3$ by this inference rule.

Second, the code $\text{Infrule}(\text{reduce_maydiff}(y), l_2)$ in Algorithm 1 outputs a command that applies the following pre-installed reduce_maydiff rule with the argument y at line l_2 :

$$\frac{(\text{reduce_maydiff}(y)) \quad y_s = e_s \quad y_t = e_t \quad \text{no reg in } e \text{ is in maydiff}}{y_s = y_t}$$

where and henceforth e_s and e_t denote the expression e with every register x replaced by x_s and x_t , respectively. This rule removes y from the maydiff set by inferring the equality $y_s = y_t$. In Figure 3, the register y is removed from the maydiff set by this rule because we have $y_s = \text{add } a_s \ 3$ and $y_t = \text{add } a_t \ 3$ with the register a not in the maydiff set.

Finally, the last step of inference at line 20 between the invariant in $[\]$ and that in $\{\}$ is trivial and performed automatically by the validator.

3. Key Ideas behind ERHL

We present key ideas for validating various translations using examples taken from our experiment. For simplicity of presentation, we will henceforth use C-like syntax instead of the LLVM syntax (e.g., $x := y + z$ for $x := \text{add } y \ z$).

3.1 Aligning Instructions with Logical No-ops

There are optimizations that break the one-to-one correspondence between the source and target instructions such as dead code elimination (DCE). For such optimizations, we align the instructions by inserting the *logical no-op* instruction lnop .

The instruction lnop is logical because it is absent from the real code and used only for validation purpose. In validation, it is interpreted as doing nothing (i.e., no-op). For

example, in the following DCE translation, we can insert lnop in the target to align the source and target instructions.

```
...                               ...
y := x + 42;  ~~~~ (lnop;)  // y not used
...                               ...
```

3.2 Validating Memory Optimizations

We show how ERHL handles non-alias information between pointers and private memory locations in order to validate memory optimizations.

Simple Alias Analysis Optimizations about memory operations are hard because pointers may be *aliased*, i.e., pointer values of different registers may point to the same location. For example, the load-load optimization in the `instcombine` pass of LLVM performs the following translation because p and q are not aliased.

```
10:  p := alloc(4);           p := alloc(4);
11:  *p := 42;                *p := 42;
12:  *q := 37;                ~~~~ *q := 37;
13:  r := *p;                 (lnop;)
14:  foo(r);                  foo(42);
```

To reason about such translations, we introduce the *noalias* predicate: $p \perp q$ denotes that the values in the registers p and q are not aliased (i.e., pointing to disjoint memory blocks) if they are pointers. We also adapt the post-invariant generator: all propositions about $*p$ are preserved by a store instruction to $*q$ if we have $p \perp q$ in the pre-invariant; otherwise all dropped. Hence the above translation can be validated in ERHL, if we can derive $p_s \perp q_s$ before $*q := 37$.

To derive and reason about *noalias* predicates, we introduce a stronger predicate, $\text{Unique}(p)$. It asserts that the value in the register p is not aliased with any values in the other registers or in the memory. Accordingly, the post-invariant generator introduces $\text{Unique}(p)$ after $p := \text{alloc}(\dots)$ and drops it after (almost) all other instructions. Also we install the following inference rules:

$$\frac{(\text{UNIQUE-NOALIAS}) \quad \text{Unique}(p) \quad p, q \text{ different registers}}{p \perp q}$$

Using such predicates, we can validate the above translation as follows. We have: after line 10, $\text{Unique}(p_s)$ as a post-invariant and $p_s \perp q_s$ by (UNIQUE-NOALIAS); after line 11, $*p_s = 42$ as a post-invariant and $p_s \perp q_s$ preserved; and after line 12, $*p_s = 42$ preserved thanks to $p_s \perp q_s$. Thus after line 13, we can infer that $r_s = 42$ using $*p_s = 42$, from which the validator can check that `foo` is invoked with 42 in both source and target.

Note that the LLVM infrastructure provides an alias analysis module that performs more complex analysis. Supporting it in ERHL is future work (see §7.2).

Private Memory Another challenge in validating memory optimizations is to handle unknown code. For example, consider the following register promotion translation performed by the *sroa* pass⁴, where the address in *p* is used locally.

```

10: p := alloc(4);      (lnop;)
11: *p := 42;           (lnop;)
12: bar();              ~~~ bar();
13: r := *p;            (lnop;)
14: foo(r);             foo(42);

```

The *sroa* pass performs this translation even though the function *bar* is unknown at compile time because the address in *p* is newly allocated and not escaped to external functions and thus the function *bar* has no way to access the address.

To make this argument formal, we introduce the *private memory* predicate: $\text{Private}(p_s)$ denotes that the register *p* contains a pointer value pointing to a *private* block in the following sense. First, in our semantic model for validation, we have the notion of *equivalent* blocks between the source and target memories, which is technically a CompCert-style memory injection (see §3.5 for more details). Then, a block in the source memory is *public* if it has a corresponding equivalent block in the target memory; otherwise *private*. Note that the values of the registers not in the maydiff set and the values stored in the public blocks are also public.

Regarding the predicate, the post-invariant generator is adapted in three ways. First, $\text{Private}(p_s)$ is introduced after the line with $p := \text{alloc}(\dots)$ in source and *lnop* in target. Second, $\text{Private}(p_s)$ is preserved by every instruction since the validator checks that only public values are passed to external functions and stored in public blocks.⁵ Third, all propositions about $*p$ are preserved by a function call if we have $\text{Private}(p_s)$ in the pre-invariant; otherwise all dropped.

Now we can validate the above translation as follows. We have: after line 10, $\text{Private}(p_s)$ as a post-invariant; after line 11, $*p_s = 42$ as a post-invariant and $\text{Private}(p_s)$ preserved; and after line 12, $*p_s = 42$ preserved thanks to $\text{Private}(p_s)$. Thus after line 13, we can infer that $r_s = 42$ using $*p_s = 42$, from which the validator can check that *foo* is invoked with 42 in both source and target.

Note that in the presence of both function calls and store instructions, we need to use all three predicates: $\text{Unique}(p)$, $p \perp q$ and $\text{Private}(p)$.

3.3 Encoding Relational Predicates using Ghosts

We show how to encode a more general form of relational invariant, $e_s = e'_t$, which is often needed, for instance, to handle loops.

The idea is to use *ghost registers* to relate the source and target expressions. Specifically, we can encode $e_s = e'_t$ as

follows: for a fresh ghost register \hat{g} ,

$$e_s = \hat{g}_s \wedge \hat{g}_t = e'_t \wedge \hat{g}_s = \hat{g}_t.$$

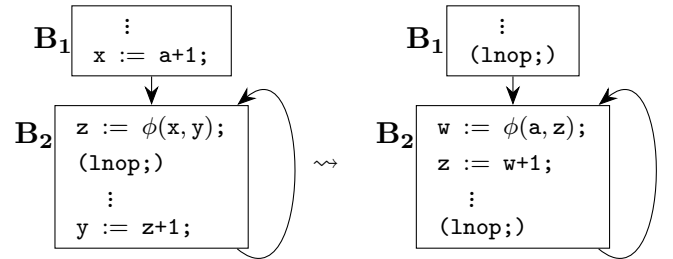
Here, $e_s = \hat{g}_s$ is a source invariant, $\hat{g}_t = e'_t$ is a target invariant and $\hat{g}_s = \hat{g}_t$ denotes that \hat{g} is not in the maydiff set.

Ghost registers are *logical* registers that do not exist in the physical program state. Thus they are existentially quantified in the semantics of ERHL invariants. Specifically, a pair of source and target states (S, T) satisfies an ERHL invariant P , if there *exists* a pair of source and target ghost register states (\hat{R}_s, \hat{R}_t) such that the pair of S extended with \hat{R}_s and T extended with \hat{R}_t satisfy P .

3.4 Handling Phinodes and Loops

We show how we handle phinodes and loops using ghost variables. Handling loops is more challenging because a register may be updated multiple times in the presence of loop even with the SSA property.

We illustrate the idea using an example. Consider the following translation performed by the *fold-phi* optimization in the *instcombine* pass of LLVM.



In the source program, the block B_1 flows into B_2 , and the block B_2 makes a self loop. The *phinode* $z := \phi(x, y)$ is a construct of SSA, whose semantics is that the register *z* is assigned the value of *x* when control coming from B_1 ; and the value of *y* when control coming from B_2 . Thus at each iteration of the loop the value of *z* gets incremented by 1.

This translation can be validated as follows. We first assert: $\text{MD}(\{x, y, w\})$ throughout B_1 and B_2 (because *x*, *y* are eliminated and *w* is introduced in the target); $x_s = a_s + 1$ at the end of B_1 ; and $y_s = z_s + 1$ at the end of B_2 . Then, after the phinode in B_2 , we assert $\text{MD}(\{x, y, w, z\})$ and $z_s = w_t + 1$, which is encoded as $z_s = \hat{t}_s + 1 \wedge \hat{t}_t = w_t$ using a fresh ghost register \hat{t} .

Now we can validate these asserted invariants. There are two non-trivial steps.

First step The invariant $z_s = \hat{t}_s + 1 \wedge \hat{t}_t = w_t$ after the phinode is validated independently for each incoming control flow.

When control coming from B_1 , before the phinode, the invariant $a_s = \hat{t}_s \wedge \hat{t}_t = a_t$ is inferred by the following rule because *a* is not in the maydiff set.

$$\begin{array}{c}
(\text{INTRO-GHOST}) \\
\frac{\hat{t} \text{ is fresh} \quad \text{no reg in } e \text{ is in maydiff}}{e_s = \hat{t}_s \wedge \hat{t}_t = e_t}
\end{array}$$

⁴ The *sroa* pass performs register promotion and replaces the *mem2reg* pass.

⁵ However it is possible to install an inference rule that eliminates $\text{Private}(p_s)$ turning a private block into a public one under certain conditions because it is allowed in our semantic model.

This rule is sound because the ghost variable \hat{t} is not used elsewhere (*i.e.*, fresh) and existentially quantified. Then $x_s = \hat{t}_s + 1$ is inferred from $x_s = a_s + 1$ from \mathbf{B}_1 and $a_s = \hat{t}_s$ by the following rule.

$$\frac{\text{(RENAME)} \quad \frac{x = e \quad y = z}{x = e[y \mapsto z]}}{x = e[y \mapsto z]}$$

After the phinode, $z_s = \hat{t}_s + 1$ is inferred from the post-invariant $z_s = x_s$ and the preserved invariant $x_s = \hat{t}_s + 1$ by transitivity. Also $\hat{t}_t = w_t$ is inferred from the preserved invariant $\hat{t}_t = a_t$ and the post-invariant $a_t = w_t$ by transitivity.

When control coming from \mathbf{B}_2 , before the phinode, the invariant $z_s = \hat{t}_s \wedge \hat{t}_t = z_t$ is inferred by (INTRO-GHOST). Then $y_s = \hat{t}_s + 1$ is inferred from $y_s = z_s + 1$ from \mathbf{B}_2 and $z_s = \hat{t}_s$ by (RENAME). After the phinode, $z_s = \hat{t}_s + 1$ is inferred from the post-invariant $z_s = y_s$ and the preserved invariant $y_s = \hat{t}_s + 1$ by transitivity. Also $\hat{t}_t = w_t$ is inferred from the preserved invariant $\hat{t}_t = z_t$ and the post-invariant $z_t = w_t$ by transitivity.

Second step Now the invariant $z_s = z_t$ (*i.e.*, removing z from the maydiff set) is validated after the second line of \mathbf{B}_2 . After $z := w + 1$, the invariant $z_t = \hat{t}_t + 1$ is inferred from the post-invariant $z_t = w_t + 1$ and the preserved pre-invariant $\hat{t}_t = w_t$ by (RENAME). Since the pre-invariant $z_s = \hat{t}_s + 1$ is also preserved, we have $z_s = z_t$ by `reduce_maydiff(z)`.

As we have seen, reasoning in the presence of phinodes is somewhat involved. In order to reduce manual effort, we have equipped the validator with an engine that automatically apply (up to some level) relevant inference rules regarding phinodes. Note that the inference rules used by the engine are parts of TCB, but the automation algorithm itself is not.

3.5 Generalizing Equality

Technically, the correctness of a translation is not equivalence but refinement between the behaviors of the source and target programs. More specifically, the behaviors of the target program need to be a subset of those of the source. Since compiler optimizations indeed exploit this definition, we also reflect it in ERHL by weakening the notion of equality.

For this, we basically follow the ideas from CompCert [13]. We use the notions of *lessdef* and *memory-injection*.

Lessdef The LLVM IR has the notion of *poison* value, which is equivalent to the *undef* value of CompCert. The poison value comes as an outcome of erroneous operation such as integer addition overflow. Since compilers have freedom to translate any erroneous behavior into an arbitrary behavior, the poison value can be translated into an arbitrary value.

Such translations may break the equality relationship between the source and target values. For example, consider

the following translation:

$$\begin{array}{ll} x := a - 1; & x := a - 1; \\ y := a - x; & \rightsquigarrow y := a - x; \\ z := y + 1; & z := 1 + 1; \end{array}$$

The register y is replaced by 1 because $y_s = a_s - x_s = a_s - (a_s - 1) = 1$ for any integer value of a_s . However, if a_s is *poison*, the equality $y_s = 1$ does not hold because arithmetic operations over *poison* unconditionally yield *poison* and thus $a_s - (a_s - 1) = \text{poison}$. As a result, we do not have $z_s = z_t$ in this valid translation because $z_s = \text{poison} \neq 2 = z_t$ when a_s is *poison*.

In order to properly capture the relationship between source and target values, we use the CompCert-style *lessdef* relation [14] instead of the equality. Specifically, a is *less defined* than b , denoted $a \sqsubseteq b$, if a is *poison* or equals to b . For example, we have $y_s = a_s - x_s = a_s - (a_s - 1) \sqsubseteq 1$ and $z_s = y_s + 1 \sqsubseteq 1 + 1 = z_t$, which justifies the above translation.

Memory Injection The memory blocks in source and target may not have one-to-one correspondence since optimizations may remove, add or merge memory allocations (*e.g.*, register promotion and allocation). Hence the corresponding pointer values in source and target may not have the same address.

In order to capture this relationship between pointer values in source and target, we also use the CompCert-style *memory injection* relation [14] instead of the equality. The idea is two-fold. First, we use logical notions of memory block and address in the operational semantics of LLVM IR. Then, in the simulation relation used to prove behavioral refinement between source and target programs, we keep a relation satisfying certain properties, called *memory injection*, between the logical memory blocks in source and target to capture the notion of *equivalent* blocks.

Though we use *lessdef* with *memory injection* in all our work, we use the equality instead in the paper for simplicity of presentation. However, all ideas presented in the paper naturally generalize to the setting with *lessdef* and *memory injection*. For example, after $x := a + b$, the post-invariant generator introduces both $x \sqsubseteq a + b$ and $a + b \sqsubseteq x$.

4. Validation of LLVM Optimizations

We show how we combine the key ideas presented in the previous section to validate translations performed by three LLVM passes: instruction combining (`instcombine`), global value numbering with partial redundancy elimination (`gvn`), and register promotion (`sroa`).

4.1 Instruction Combining

The `instcombine` pass in LLVM is a collection of many micro-optimizations that detects inefficient patterns of instructions and transforms them into a cheaper form. It is one of the biggest passes in LLVM, comprising 23,000 lines of

C++ code with more than 1000 micro-optimizations [15]. Among those micro-optimizations, we extend and apply our validator to selected 146 micro-optimizations. We select at least one micro-optimization from each of the 10 files in the instcombine pass.⁶

The selected 146 micro-optimizations are not just very simple “peephole” optimizations. They include non-trivial optimizations: those across phinodes and loops such as fold-phi and memory optimizations such as load-load, load-store, and dead-store-elim. Indeed all the ideas of §3 are used as follows. Out of the 146 micro-optimizations:

- 30 use the logical no-op instruction (§3.1);
- 3 use the noalias predicate (§3.2);
- 4 use relational predicates with ghost registers (§3.3, §3.4);
- 5 materialize poison into a concrete value (§3.5).

In the supplementary material [1, §C], we present the full list of the 146 micro-optimizations with the techniques used by them.

4.2 Global Value Numbering with Partial Redundancy Elimination

Global Value Numbering The gvn pass in LLVM performs global value numbering (GVN), which eliminates redundant instructions whose outcome is statically known equivalent to that of a previous instruction. The GVN algorithm identifies redundant instructions by classifying expressions into equivalence classes of those yielding the same outcome.

As an example, consider the following translation:

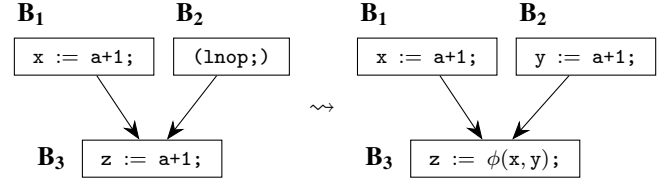
```
x := a + 1;      x := a + 1;
y := a + 1;      (lnop;) // y replaced by x
z := y + 2;      ~~~ z := x + 2;
v := x + 2;      (lnop;) // v replaced by z
w := z + v;      w := z + z;
```

Here, x and y are classified as results of equivalent expressions and thus y is removed and replaced by x . Subsequently z and v are classified as equivalent and thus v is removed and replaced by z .

The validation of GVN is similar to that of simple arithmetic optimizations like add-assoc, but it is required to generate proofs in a recursive fashion, due to the recursive nature of GVN’s value classification algorithm.

Partial Redundancy Elimination The gvn pass also performs *partial redundancy elimination* (PRE), which eliminates partial redundancy using the classification result of GVN.

As an example, consider the following translation:



Here, z in B_3 is partially redundant because it is redundant only when the preceding block is B_1 . PRE removes such partial redundancy as follows. First, it inserts $y := a+1$ in B_2 to make z fully redundant and then removes the redundancy by transforming $z := a+1$ into $\phi(x, y)$.

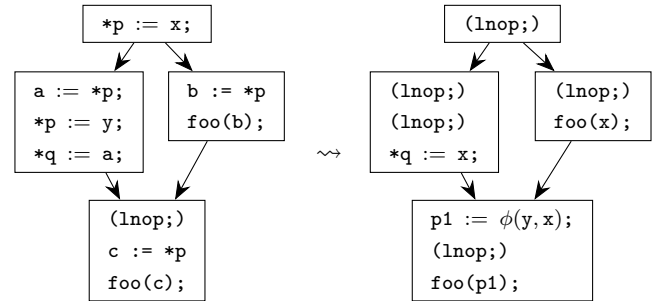
For validation of PRE, we additionally use the idea of §3.3 and §3.4 to perform validation across phinodes and loops. Note that PRE is performed even across loops.

In our experiment, we fully cover the main algorithm, GVN and PRE, of the gvn pass, which consists of 1,496 lines of C++ code. However, the gvn pass also performs several micro-optimizations that help the main GVN-PRE algorithm, which we do not currently handle.

4.3 Register Promotion

The register promotion optimization optimizes loads and stores to memory locations into reads and writes to registers, provided that the locations are used only locally (*i.e.*, their addresses are not escaped to external functions nor to the memory). Furthermore, it is this register promotion that essentially performs the SSA transformation, since stores to the same location are allowed statically multiple times whereas writes to a register is allowed statically once. This algorithm is implemented in a single file consisting of 989 lines of C++ code, which is invoked by the sroa pass (previously, invoked by the mem2reg pass).

As an example, consider the following translation in which the location p is promoted to existing registers x , y and a fresh register $p1$:



To validate this translation, most of the ideas in §3 are used. First, the logical no-op instruction (§3.2) are properly inserted. Second, the noalias and private memory predicates (§3.2) are used for preserving invariants about $*p_s$ after the store instruction $*q := a$ and the function calls $foo(b)$ and $foo(c)$. Finally, ghost registers (§3.3, §3.4) are used to express relational predicates such as $*p_s = p1_t$ after the phinode $p1 := \phi(y, x)$. Such phinodes are inserted in the

⁶ We omitted the files on vectors, function calls, and bit operations.

presence of branches and loops to keep the target program in SSA form.

Note that the implementation of register promotion in LLVM temporarily breaks the SSA and semantics preservation properties, which makes step-by-step verification hard [25]. In order to sidestep this problem, we validate the whole steps of a register promotion translation at once. More specifically, we accumulate all proofs generated during the whole translation and merge them as a single proof, which we pass to the validator.

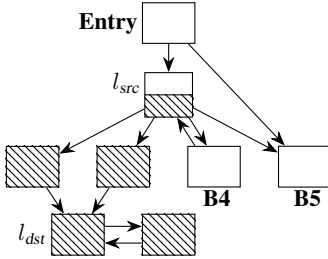
5. Validator Components

In this section, we present further details of the components of our ERHL system.

5.1 Invariant Propagator

As discussed in §2.2.1, the invariant propagator generates full invariants from commands given by the compiler. For example, the command `Propagate(src, x = add y 42, lsrc, ldst)` requires to add the proposition $x = \text{add } y \ 42$, say P shortly, to the source invariants between l_{src} and l_{dst} , where the notion of program point (e.g., l_{src} , l_{dst}) is generalized from a line number to a block name with a line number in it.

More specifically, the proposition P should be added at every program point appearing in a path from l_{src} to l_{dst} that does not visit l_{src} but may visit l_{dst} in-between. Since l_{src} is the source of the proposition P so that we can get P as a post-condition every time we visit l_{src} , there is no need to propagate P along a path from l_{src} to l_{src} . For example, consider the following program.



The marked area between l_{src} and l_{dst} is where we should propagate the proposition P .

Thanks to the SSA property [9], we can efficiently calculate the program points between l_{src} and l_{dst} . First, we can assume that l_{src} dominates l_{dst} (i.e., one should have visited l_{src} to reach l_{dst} from the entry point), since the proposition P created at l_{src} should hold at l_{dst} . Then a program point l is on a path from l_{src} to l_{dst} that does not visit l_{src} in-between if and only if (i) l_{src} dominates l and (ii) l_{dst} is reachable from l without visiting l_{src} . We efficiently check the first condition using the dominator tree [9] and the second condition by a backward BFS search from l_{dst} .

In the above example, any program point in the marked area is such that l_{src} dominates it and l_{dst} is reachable from it without visiting l_{src} . On the other hand, l_{dst} is not reachable

from the block B_4 without visiting l_{src} , and l_{src} does not dominate the block B_5 .

Note that the invariant propagator is not a part of TCB: validation may fail but cannot succeed incorrectly due to bugs in the propagator.

5.2 Post-Invariant Generator

The post-invariant generator, `postinv : Instr × Instr × Inv → option Inv`, computes a post-invariant of the given pre-invariant after the given source and target instructions. As illustrated in §2.2.2, we define it as follows:

```

postinv(is, it, P) =
  if check(is, it, P)
  then Some auto(add(is, it, remove(is, it, P)))
  else None .

```

The function `check(is, it, P)` checks whether the instructions i_s and i_t produce the same observable events under the program states satisfying the pre-invariant P . Specifically, it checks whether i_s and i_t call the same function with equivalent arguments, if one of them is a function call; and whether i_s and i_t store equivalent values to equivalent memory locations including global variables, if one of them is a store to a public location.

The function `check` checks the equivalence of a register x in source and y in target by checking whether there is a (possibly ghost) register z such that $x_s = z_s \wedge y_t = z_t$ and z is not in the maydiff set. This is enough because we can always make provably equivalent registers in source and target to factor through a fresh ghost register by transitivity and (INTRO-GHOST).

The function `remove(is, it, P)` computes a post-invariant Q by removing from the pre-invariant P those propositions that are invalidated by executing the instructions i_s and i_t . We have mostly shown what invariants are invalidated in §2 and §3. Here we explain one important case that is not clearly presented so far. If i_s is $x := e$, all propositions about x_s are removed; and furthermore x is unconditionally added to the maydiff set, which may be dropped later by `auto` in case i_t executes an equivalent instruction. We proceed similarly for i_t .

The function `add(is, it, Q)` strengthen the given post-invariant Q by adding new propositions obtained from the fact that i_s and i_t are just executed. We have mostly shown what invariants are newly added in §2 and §3. Here we explain one important case that is not clearly presented so far. If both i_s and i_t execute $p := \text{alloc}(\dots)$ with equivalent arguments, then the register p is removed from the maydiff set.

Finally, the function `auto(Q)` automatically applies pre-defined sequences of inference rules to the given post-invariant Q to get a stronger post-invariant. First, `auto` in our validator applies inference rules regarding phinodes and loops, as discussed in §3.4. Second, `auto` applies the `reduce_maydiff` rule to the register just defined in order to

	Compiler (Covered)	Proof Generation
instcombine	799	4950
gvn	1496	1207
sroa	989	1872
Library	-	17020 (auto-gen'd: 11007)

Figure 4. Lines of Compiler and Proof Generation Code

remove it from the maydiff set in case the same instruction is executed in source and target.

For example, suppose that i_s and i_t are $x := y + 1$ and y is not in the maydiff set of the pre-invariant. Even so, the function `add` unconditionally adds x to the maydiff set of the post-invariant. However, the function `auto` removes the register x from the maydiff set by applying `reduce_maydiff(x)`.

Note that the functions `check`, `remove` and `add` are parts of TCB. For the function `auto`, only the inference rules used by `auto` are parts of TCB, but the automation algorithm of `auto` itself is not.

5.3 Inference Rules

The function `apply` : $InfRule \times Inv \rightarrow Inv$ defines how the inference rules work. The set $InfRule$ consists of inference rule names with arguments and the function `apply` is defined separately for each inference rule. In order to extend our validator with new inference rules, we just need to add more rules in the set $InfRule$ and add to `apply` the definitions (*i.e.*, functions from Inv to Inv) of the newly added rules.

The `apply` function is a part of TCB and should be sound in the sense that for any inference rule r , any source and target program states satisfying an invariant I should also satisfy the invariant `apply(r, I)`.

6. Results

6.1 Development

Our validator targets on 146 micro-optimizations in the `instcombine` pass, the main GVN-PRE algorithm in the `gvn` pass, and the register promotion algorithm in the `sroa` pass in LLVM 3.7.1. We just added the proof generation code to the original compiler without changing any existing code.

Figure 4 shows the number of lines of the compiler code in C++ for which we inserted the proof generation code. We separately counted the number of lines of our library for writing proof generation code. 64.7% of the library is about printing commands for inference rules and are automatically generated from 1,835 lines of code in a custom DSL.

The proof generation code is long compared to the covered compiler code, and in particular, it is as twice as long for register promotion in `sroa`. This is mainly because we needed to reimplement the compiler’s logic in the proof generation code in order to keep the original compiler code untouched. Another reason is that the proof generation code contains many boilerplate code. We expect it can be reduced significantly.

The development of our validator took approximately 30 person-months in total. It took 5 months to upgrade and stabilize Vellvm, such as upgrading memory model and fixing a memory leak. It took 1 month to implement the validator, and 3 months to formally verify it. It took 2 months to implement the C++ library for proof generation. It took 2 month to prepare test infrastructure for the validator. It took 4 months to validate `instcombine`, 4 months to validate `gvn`, and 9 months to validate `sroa`. It took longer to validate optimization passes than expected, mainly because the compiler code is not designed in consideration of verification.

6.2 Performance Evaluation

We evaluate the performance of our validator in our experiments, in which we validated the compilation of the SPEC CINT2006 C Benchmarks [5] and other 5 open-source C projects (the biggest benchmarks used in [17]⁷), totalling 3.9 million lines of C code. We omitted 4 files from the benchmarks: 3 of them contain instructions currently not supported in Vellvm, including the `indirectbr` instruction, and one is simply too big (151MB) for our validator to handle. We conducted the experiments in a Linux workstation with Intel Xeon E5-2630 CPU (2.6GHz, 12 cores with hyper-threading), 128GB RAM, and 1TB SSD (Samsung 850 PRO).

Figure 5 summarizes the experimental results for each optimization pass. In the experiment, we compiled each benchmark program using the original and the modified compilers with the `-O2` option. We counted the number of the generated validation proofs (**#V**), those proofs our validator failed to validate (**#F**), and those currently not supported (**#NS**). Everything else (**#V** - **#F** - **#NS**) is successfully validated. We also measured the time spent to perform each optimization in the original compiler (**Comp**), to perform each optimization and calculate validation proofs in the modified compiler (**PCal**), to print the source and the target programs and the proofs to files (**Print**), and to validate the generated proofs by the validator (**Vali**). When measuring the times, we maximized the CPU utilization by executing 24 jobs in parallel (because the compiler and the validator are single-threaded programs), measured the CPU time (User+System) used by each job, and summed up them to get the total.

Out of 1,191,063 validations in total, 993,157 (83.4%) are successfully validated. All 254 (0.02%) failures are due to compiler bugs: 247 are due to two bugs in `gvn`, and 7 are due to a bug in `sroa`. We present the three compiler bugs in the supplementary material [1, §A].⁸ The other 197,652 (16.6%) translations are currently not supported in our validator, among which 191,819 (97.0%) use instructions not supported in Vellvm: vector operations (181,640, 91.9%), debug attributes (8,485, 4.3%) and atomic fences (1,694, 0.9%); 2838 (1.4%) use the alias analysis module of LLVM;

⁷ We omitted Linux, since it is currently not compiled with LLVM. See [3] for more details.

⁸ We do not identify the bugs for the double-blind review.

	LOC	Result									Time (sec.)											
		instcombine			gvn			sroa			instcombine				gvn				sroa			
		#V	#F	#NS	#V	#F	#NS	#V	#F	#NS	Comp	PCal	Print	Vali	Comp	PCal	Print	Vali	Comp	PCal	Print	Vali
400.perlbench	168.16K	59753	0	7357	11882	17	0	1745	0	1	5.78	36.52	4063.25	374.44K	1.39	2.27	144.83	142.51K	0.24	54.07	123.17	96.66K
401.bzip2	8.29K	4545	0	1865	1681	0	0	90	0	0	0.36	4.16	137.19	2.82K	0.08	0.23	7.30	2.60K	0.01	3.48	5.99	34.98K
403.gcc	517.52K	140730	0	5076	36999	20	13	5263	0	5	18.59	156.44	15755.35	657.53K	5.40	7.87	365.96	124.51K	0.71	153.84	324.60	446.02K
429.mcf	2.69K	492	0	57	149	0	0	24	0	0	0.05	0.13	1.53	0.02K	0.01	0.02	0.16	0.01K	< 0.01	0.29	0.50	0.02K
433.milc	15.04K	3840	0	517	2046	0	0	235	0	2	0.46	1.02	34.77	0.31K	0.10	0.14	4.30	0.19K	0.03	2.27	5.27	0.35K
445.gobmk	196.24K	16324	0	1724	7176	0	56	2641	0	1	2.99	8.27	1325.91	12.75K	0.70	0.95	34.65	6.31K	0.18	35.88	133.45	51.61K
456.hammer	35.99K	11123	0	3433	3538	3	3	558	0	0	1.48	4.05	156.60	1.10K	0.31	0.39	7.42	0.55K	0.07	10.39	21.85	2.99K
458.sjeng	13.85K	3885	0	297	1751	0	1	130	0	0	0.63	1.65	71.29	1.14K	0.15	0.17	4.57	0.44K	0.02	3.01	5.62	1.03K
462.libquantum	4.36K	1052	0	793	357	0	265	123	0	79	0.15	0.34	7.39	0.04K	0.03	0.04	0.52	0.01K	0.01	0.68	1.41	0.01K
464.h264ref	51.58K	22885	0	5529	12858	27	4	532	1	0	2.78	21.15	1306.39	16.31K	0.71	1.54	104.64	14.08K	0.10	26.96	59.29	88.88K
470.lbm	1.16K	174	0	51	77	0	0	19	0	0	0.03	0.07	1.38	0.01K	0.01	0.01	0.07	< 0.01K	< 0.01	0.37	0.67	0.54K
482.sphinx3	25.09K	6280	0	1056	1652	0	3	364	0	0	0.84	1.93	62.05	0.56K	0.15	0.18	3.13	0.23K	0.04	3.99	8.41	0.69K
sendmail-8.15.2	138.68K	14172	0	411	4647	0	120	536	0	403	2.24	7.30	737.00	11.52K	0.52	0.65	37.55	5.47K	0.14	8.32	20.94	0.40K
emacs-25.1	463.54K	110184	0	9605	28730	15	35	5149	0	4	15.16	63.08	11014.64	128.56K	2.77	4.67	395.55	48.98K	0.65	149.73	346.02	232.21K
python-3.4.1	486.38K	89612	0	13419	28014	117	51	8784	0	89	15.97	54.27	6998.85	51.02K	3.75	5.24	284.35	27.81K	0.85	92.95	260.57	40.95K
gimp-2.8.18	1004.20K	151905	0	44676	38282	22	438	19448	6	528	29.03	80.44	3455.24	19.94K	4.99	7.29	153.14	8.57K	1.58	125.91	362.10	29.00K
ghostscript-9.14.0	797.65K	247471	0	83540	68166	26	6969	12996	0	9176	27.75	95.44	8502.72	74.12K	7.59	11.17	475.71	40.29K	2.08	58.17	176.28	1.50K
Total	3930.42K	884427	0	179406	248005	247	7958	58637	7	10288	124.29	536.25	53631.53	1352.19 K	28.67	42.81	2023.85	422.56K	6.75	730.31	1856.16	1027.83K

Figure 5. Experimental Results

2200 (1.1%) alter type declarations; 693 (0.4%) require deeper analysis on functions such as read-only function analysis; and 102 (< 0.1%) are beyond the scope of the main GVN-PRE algorithm.

Proof generation and validation are slow compared to compilation. In particular, on average, proof calculation took 4.3 times for instcombine, 1.5 times for gvn, and 108.2 times for sroa, compared to the original compilation time.

Since our proof generation code does not have bigger time complexity than the optimization algorithm itself has, we believe we can significantly reduce the time for generating and validating proofs. There is plenty of room for improvement in calculating proofs. For example, the proof calculation code creates unnecessarily many C++ objects and auxiliary data. In particular, sroa calculates too verbose proofs (the average size of a proof is 343.42KB), mainly because we validated the whole step of a register promotion translation at once. There are opportunities of splitting the whole proof into many smaller ones avoiding the problem of broken SSA (see §4.3), which we believe will reduce the proof generation time significantly. Furthermore, we can pass programs and proofs from the compiler to the validator faster, (i) by compressing the data by printing only difference between the source and the target programs and making the proofs more terse, and (ii) by writing data to RAM instead of SSD. Also, the validator is naively implemented in order to make verification of the validator easier. Finally, we can massively parallelize the validations in the clouds, as each validation job is independent from each other.

6.3 Verification

In order to reduce our validator’s TCB, we formally verified major parts of our validator in the Coq proof assistant.

Semantics We use the formal semantics of LLVM IR from the Vellvm project [26]. Vellvm used the CompCert memory model [14] version 1.9 and we upgraded it to version 2.4 in order to use the notion of *permission* in the LLVM semantics. We also added the *switch* instruction to the formalization

of LLVM IR, which is used by 50.8% of our benchmark C projects used for performance evaluation.

Vellvm still lacks the support for various features in LLVM, including vector operations, atomic operations for concurrency, and attributes such as *noalias* and *readonly*. Furthermore, casts between integers and pointers are not properly formalized in Vellvm and the CompCert memory model. However, the idea of Kang *et al.* [11] to model integer-pointer casts is applicable to our memory model, which we leave as future work.

Finally, the *undef* value, a weaker notion than *poison*, of LLVM IR is simply formalized as *poison* in Vellvm. In fact, it is well-known in the LLVM community that the semantics of *poison* together with *undef* is seriously broken [4]. Due to the lack of a weaker notion of *undef* semantics, some optimizations, such as loop unswitching, are not valid in our model as well as in CompCert. Further research is needed to solve this problem.

Verification of Validator We proved in Coq the correctness of the post-condition generator of our validator using simulation relations similar to those used in CompCert, currently leaving many admits on minor details. The main theorem states that if our validator validates a proof of a translation, the source and target programs are related by a *simulation* relation. We also proved the adequacy of simulation relations (*i.e.*, simulation implies behavioral refinement). Our correctness result is the same as that of CompCert.

Inference Rules We installed 194 inference rules for our experiment. None of them are verified but we believe they are clearly sound because they are all simple rules. In fact, we proved about 30 inference rules sound in Coq for an early version of our development, which we have not ported to the current version yet.

7. Discussion

7.1 Related Work

There has been many proposals for compilation verification/validation, but our work is the first verified approach to mainstream compilers.

(Foundational) proof carrying code (PCC) [16, 6] is a (verified) proof system for validating safety properties of a program. We are inspired by PCC’s use of proofs for validating program properties, but we validated the relational property of semantic preservation between source and target programs, rather than unary properties of a single program.

Tristan *et al.* [20] and Stepp *et al.* [19] developed translation validators for LLVM optimization passes, including dead code elimination, GVN-PRE, constant propagation, and loop invariant code motion. We are also inspired a lot by these approaches. However, these differ from our approach in several aspects. First, their validators are inherently incomplete, since they are not given any extra information from compilers. Hence they failed to validate about 20% of programs in the benchmark. Furthermore, they do not cover register promotion, which requires complex memory reasoning, nor formally verify their validators.

CompCert [13] is the first formally verified optimizing C compiler. However, CompCert implements a relatively small number of optimizations compared to mainstream compilers such as LLVM. Also, CompCert is implemented in Coq [2] and run in OCaml, which is a purely functional language, and thus its compilation time is relatively slow compared to the mainstream compilers written in C++. We think the verified compiler approach is not very adequate for mainstream compilers because the verification cost is high and the compilation time is slow.

Zhao *et al.* [26, 25] implemented and verified the `vmem2reg` pass in Coq for the Vellvm project. The `vmem2reg` pass performs register promotion, but its algorithm is significantly simplified compared to that in LLVM.

Translation validators are verified for various optimizations in CompCert. Tristan and Leroy verified ones for instruction scheduling [21], lazy code motion [23], software pipelining [22]; Rideau and Leroy, register allocation [18]; Barthe *et al.* [7], SSA transformation; and Demange *et al.* [10], GVN and sparse conditional constant propagation (SCCP). These validators provide the same level of correctness guarantee as verified compilation. Compared to our approach, the difference is that these validators are highly specialized for the target optimizations, without sharing any common validator like our ERHL validator.

Lopes *et al.* [15] presented Alive, a domain-specific language for writing peephole optimizations. Using Alive, one can either prove the correctness of an optimization algorithm using SMT solvers, or find a counterexample. The authors ported 300 micro-optimizations in `instcombine` to Alive, and in doing so they found 8 bugs in `instcombine`. While practical, Alive is limited to supporting only peephole op-

timizations that are performed inside a single block (*i.e.*, programs with linear control flow). Furthermore, Alive relies on SMT solvers and makes simplifying assumptions on the LLVM semantics, and thus provides a lower level of correctness guarantee than that of compilation verification performed in Coq.

CSmith [24, 8] and EMI [12] are random testing systems for validating C compilers. They improved the reliability of mainstream C compilers by finding hundreds of bugs in them. However, testing does not assure the absence of bugs.

7.2 Future Work

There are a number of interesting issues remaining for future work.

Changing Control-Flow Graph So far we assumed that the source and the target programs can be aligned by inserting logical no-op instructions. While this assumption applies to the majority of LLVM optimization passes (see the supplementary material [1, §B] for more details), some important optimizations, including loop unrolling and loop unswitching, are beyond the reach of the current paper. We believe we can lift the assumption by generally mapping the source and the target program points for the validator.

Supporting Analysis Passes Currently our validator does not support analysis passes, including alias analysis, read-only function analysis, and memory dependence analysis. We believe we can support them by adding appropriate predicates and inference rules in our validator.

Supporting Missing Features As discussed in §6.3, it is an interesting future work to support missing features in Vellvm.

Validating Frontend & Backend We focused on the IR-to-IR optimization passes, leaving the frontend (C to IR) and backend (IR to assembly) passes as future work.

8. Conclusion

In this paper, we propose a practical approach to compilation verification that applies to mainstream compilers without sacrificing the compilation performance.

The post validation does not need to be performed every time we compile because it consumes huge computing resources. We think it would most useful to perform post validation only before deploying the software or sometimes in debugging to make sure the absence of compiler bugs.

Also we believe that it would save a lot of efforts if we develop a compiler together with the validator because we can design a compiler in a way that validation become easier.

References

- [1] Anonymized development and supplementary material.
- [2] The coq proof assistant. <https://coq.inria.fr/>. Accessed: 2016-11-15.

- [3] LLVM Linux. <http://llvm.linuxfoundation.org>. Accessed: 2016-11-15.
- [4] RFC: Killing undef and spreading poison. <http://lists.llvm.org/pipermail/llvm-dev/2016-October/106182.html>. Accessed: 2016-11-15.
- [5] The spec cint2006 benchmark. <https://www.spec.org/cpu2006/CINT2006/>. Accessed: 2016-11-15.
- [6] Andrew W. Appel. Foundational proof-carrying code. LICS '01.
- [7] Gilles Barthe, Delphine Demange, and David Pichardie. Formal verification of an SSA-based middle-end for CompCert. *ACM Trans. Program. Lang. Syst.*, 36(1).
- [8] Yang Chen, Alex Groce, Chaoqiang Zhang, Weng-Keen Wong, Xiaoli Fern, Eric Eide, and John Regehr. Taming compiler fuzzers. PLDI '13.
- [9] Ron Cytron, Jeanne Ferrante, Barry K. Rosen, Mark N. Wegman, and F. Kenneth Zadeck. Efficiently computing static single assignment form and the control dependence graph. *ACM Trans. Program. Lang. Syst.*, 13(4).
- [10] Delphine Demange, David Pichardie, and Léo Stefanesco. Verifying fast and sparse SSA-based optimizations in Coq. CC '16.
- [11] Jeehoon Kang, Chung-Kil Hur, William Mansky, Dmitri Garbuzov, Steve Zdancewic, and Viktor Vafeiadis. A formal C memory model supporting integer-pointer casts. PLDI '15.
- [12] Vu Le, Mehrdad Afshari, and Zhendong Su. Compiler validation via equivalence modulo inputs. PLDI '14.
- [13] Xavier Leroy. Formal certification of a compiler back-end or: Programming a compiler with a proof assistant. POPL '06.
- [14] Xavier Leroy, Andrew W. Appel, Sandrine Blazy, and Gordon Stewart. The CompCert memory model, version 2. Research report RR-7987, INRIA, 2012.
- [15] Nuno P. Lopes, David Menendez, Santosh Nagarakatte, and John Regehr. Provably correct peephole optimizations with alive. PLDI '15.
- [16] George C. Necula. Proof-carrying code. POPL '97.
- [17] Hakjoo Oh, Kihong Heo, Wonchan Lee, Woosuk Lee, Daejun Park, Jeehoon Kang, and Kwangkeun Yi. Global sparse analysis framework. *ACM Trans. Program. Lang. Syst.*, 36(3).
- [18] Silvain Rideau and Xavier Leroy. Validating Register Allocation and Spilling. CC '10.
- [19] Michael Stepp, Ross Tate, and Sorin Lerner. Equality-based translation validator for llvm. CAV'11.
- [20] Jean-Baptiste Tristan, Paul Govereau, and Greg Morrisett. Evaluating value-graph translation validation for llvm. PLDI '11.
- [21] Jean-Baptiste Tristan and Xavier Leroy. Formal verification of translation validators: A case study on instruction scheduling optimizations. POPL '08.
- [22] Jean-Baptiste Tristan and Xavier Leroy. A simple, verified validator for software pipelining. POPL '10.
- [23] Jean-Baptiste Tristan and Xavier Leroy. Verified validation of lazy code motion. PLDI '09.
- [24] Xuejun Yang, Yang Chen, Eric Eide, and John Regehr. Finding and understanding bugs in C compilers. PLDI '11.
- [25] Jianzhou Zhao, Santosh Nagarakatte, Milo M.K. Martin, and Steve Zdancewic. Formal verification of SSA-based optimizations for LLVM. PLDI '13.
- [26] Jianzhou Zhao, Santosh Nagarakatte, Milo M.K. Martin, and Steve Zdancewic. Formalizing the LLVM intermediate representation for verified program transformations. POPL '12.