Gilles Barthe, IMDEA Software Institute Delphine Demange, IRISA - University of Rennes 1 / Inria David Pichardie, IRISA - ENS Rennes / Inria

CompCert is a formally verified compiler that generates compact and efficient code for a large subset of the C language. However, CompCert foregoes using SSA, an intermediate representation employed by many compilers that enables writing simpler, faster optimizers. In fact, it has remained an open problem to verify formally an SSA-based compiler. We report on a formally verified, SSA-based, middle-end for CompCert. In addition to providing a formally verified SSA-based middle-end, we address two problems raised by Leroy in 2009: giving an intuitive formal semantics to SSA, and leveraging its global properties to reason locally about program optimizations.

Categories and Subject Descriptors: D.3.4 [**PROGRAMMING LANGUAGES**]: Processors – Compilers; F.3.1 [**LOGICS AND MEANINGS OF PROGRAMS**]: Specifying and Verifying and Reasoning about Programs – Mechanical Verification

General Terms: Languages, Reliability, Verification

Additional Key Words and Phrases: Single Static Assignement, Compiler Verification, Mechanized Proof

ACM Reference Format:

Gilles Barthe and Delphine Demange and David Pichardie. 2014. Formal Verification of an SSA-based Middle-end for CompCert. *ACM Trans. Program. Lang. Syst.* Vol, Num, Article Art (May Year), 37 pages. DOI:http://dx.doi.org/10.1145/000000000000

Contents

1	Introduction	3					
2	Background						
	2.1 Static Single Assignment form	5					
	2.1.1 Converting into SSA form	5					
	2.1.2 Maximal, minimal and pruned SSA	6					
	2.1.3 SSA-based optimizations	7					
	2.2 CompCert	7					
3	The RTL language	8					
	3.1 Syntax and semantics	8					
	3.2 Normalizing RTL syntax	10					
4	The SSA language 4.1 SSA programs	10 10					

Partially funded by Spanish project TIN2009-14599 DESAFIOS 10, and Madrid Regional project S2009TIC-1465 PROMETIDOS, and French project ANR Verasco, FNRAE ASCERT and Bretagne Regional project CertLogS.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© Year ACM 0164-0925/Year/05-ARTArt \$15.00 DOI:http://dx.doi.org/10.1145/0000000.0000000

	4.9	4.1.1 Syntax 4.1.2 Strict SSA 4.1.3 Well-formed SSA programs Sementias	$10 \\ 11 \\ 12 \\ 13$					
	4.2	4.2.1 Exploiting normalization for an intuitive semantics 4.2.2 Parallel execution of ϕ -blocks	13 13 13					
5	Tra	nslation validation of SSA generation	14					
	5.1	Type system	15					
		5.1.1 Liveness	15					
		5.1.2 Typing rules for instructions	15					
		5.1.3 Typing rules for edges and functions	16					
	5.2	Strictness	18					
	5.3	Soundness	18					
		5.3.1 Simulation relation	19					
		5.3.2 Proof sketch	20					
	5.4	Completeness of the type system	20					
		5.4.1 Specification of Cytron <i>et al.</i> 's algorithm	21					
		5.4.2 Building a witness global typing	22					
		5.4.3 The witness global typing is a correct typing	22					
	5.5	Implementation	24					
6	The	SSA equation lomma	95					
U	6 1	Fountion lomma	20 95					
	6.2	Application with Copy Propagation	25 26					
7	Vali	idation of Global Value Numbering	26					
8	Con	oversion out of SSA	28					
	8.1	Critical edges	28					
	8.2	The swap problem	29					
	8.3	Correctness proof	29					
•	т		00					
9	Imp	Continue of the second se	49					
	9.1	Coding effort, and ease of proof	29					
	9.2		30					
		9.2.1 Efficiency of SSA validator	30					
		9.2.2 Effectiveness of GVN optimizer	31					
		9.2.3 Efficiency of the Generated code	31					
10 Related work								
	10.1	Machine-checked formalizations	32					
	10.2	2 Translation validation and type systems	34					
11	Con	clusion and future work	34					

1. INTRODUCTION

Static single assignment. Static single assignment (SSA) form [Cytron et al. 1991] is an intermediate representation where variables are statically assigned exactly once. Thanks to the considerable strength of this property, the SSA form simplifies the definition of many optimizations, and improves their efficiency, as well as the quality of their results. It is therefore not surprising that many modern compilers, including GCC and LLVM, rely heavily on SSA form, and that there is a vast body of work on SSA. However, the simplicity of SSA form is deceptive, and designing a correct SSAbased middle-end compiler has been fraught with difficulties. In fact, it has been a significant challenge to design efficient, semantics-preserving, algorithms for converting programs into SSA form, or optimizing SSA programs, or even transforming programs out of SSA form.

Verified Compilers. Compiler correctness aims to provide rigorous proofs that compilers preserve the behavior of programs they compile. After 40 years of rich history, the field is entering into a new era, with the advent of realistic and mechanically verified compilers. This new generation of compilers was initiated with CompCert [Leroy 2009], a compiler that is programmed and verified in the Coq proof assistant and generates compact and efficient assembly code for a large fragment of the C language. Leroy's CompCert has been rightfully acclaimed as a *tour de force*, but it foregoes relying on an SSA-based middle end. Leroy [2009] reports:

Since the beginning of CompCert we have been considering using SSA-based intermediate languages, but were held off by two difficulties. First, the dynamic semantics for SSA is not obvious to formalize. Second, the SSA property is global to the code of a whole function and not straightforward to exploit locally within proofs.

and adds: "A typical SSA-based optimization that interests us is global value numbering". However verifying GVN is a significant challenge, and its formal verification has remained beyond the current state-of-the-art in certified compilers.

The structural properties of SSA are well-identified in the literature, and some proofs of SSA-based analyses and transformations can be found [Cytron et al. 1991; Chow et al. 1997; Boissinot et al. 2008]. What is missing in those works is the semantic counterparts of those properties. The proofs are traditionally based on how the SSA-based algorithms work and the information they compute (i.e. properties of the CFG). In particular, the semantic properties and invariants established by the SSA generation algorithm are never expressed precisely. This is probably due to the lack of a definition of a semantics for SSA that would be both formal and close to the intuitive definition given in the seminal papers [Alpern et al. 1988; Cytron et al. 1991].

Static Single Assignment meets verified compilers. The thesis of our work is that a compiler can be realistic, verified and still rely on a SSA form. To support our thesis, we provide the first verified SSA-based middle-end. Rather than programming and proving a verified compiler from scratch, we have programmed and verified a SSA-based middle-end compiler that can be plugged into CompCert at the level of RTL. Figure 1 describes the overall architecture. Our middle-end performs four phases: (i) normalization of RTL program; (ii) transformation from RTL form into SSA form; (iii) optimization of programs in SSA form, including Global Value Numbering (GVN) [Alpern et al. 1988]; (iv) transformation of programs from SSA form to RTL form; and relies on CompCert for the transformation from C to RTL programs prior to SSA conversion, and from RTL programs to assembly code after conversion out of SSA—our goal is to develop a real-

Gilles Barthe et al.



Fig. 1: The SSA Middle-end

istic and verified SSA-based middle-end, rather than to demonstrate that SSA-based optimizations dramatically improve the efficiency of generated code.

We validate our compiler middle-end with a mix of techniques directly inherited from CompCert. We employ translation validation [Samet 1975; Pnueli et al. 1998; Necula 2000]— a technique increasingly favored by CompCert [Tristan and Leroy 2009; 2010]— for converting programs into SSA and for GVN. Specifically, we program in Coq verified checkers that validate *a posteriori* results of untrusted computations, and we implement in OCaml efficient algorithms for these computations. We rely on Cytron et al. [1991]'s algorithm for computing minimal SSA form, and on Alpern *et al.* iteration strategy [Alpern et al. 1988] for computing a numbering in GVN. In contrast, the normalization of RTL programs, and the conversion out of SSA are directly programmed and proved in Coq. In addition, our work addresses the two issues raised by Leroy [2009]. First, we give a simple and intuitive operational semantics for SSA, that follows the informal description given in [Cytron et al. 1991], and does not require any artificial state instrumentation. Second, we define on SSA programs two global properties, called strictness and equational form, allowing to conclude reasonably directly that the substitutions performed by GVN and other optimizations are sound.

Summarizing, our work provides the first verified SSA-based middle-end, the first formal proof of an SSA-based optimization, as well as an intuitive semantics for SSA. It thus serves as a good starting point for further studies of verified and realistic SSAbased compilers.

This paper supersedes [Barthe et al. 2012]. The main differences are a proof of completeness for the SSA translation validator, a description of the type inference implementation, a more detailed description of the conversion out of SSA and more precise measurements of the GVN optimizer efficacy. The companion Coq development is available online [CompCertSSA 2012].

Contents. The paper is organized as follows. Section 2 provides a brief primer on SSA and CompCert. In Section 3, we recall the syntax and semantics of RTL, the Comp-Cert intermediate representation into which we plug our middle-end, and we explain the pre-processing of RTL prior to the middle-end. Section 4 defines the SSA language used by our middle-end. Conversion to and out of SSA forms are presented in Section 5 and 8 respectively. Section 7 presents the GVN optimizer. The correctness of this optimization relies on a core lemma, which we present in Section 6. We conclude with experimental results in Section 9 and related work in Section 10. Throughout the paper, we use Coq syntax for our definitions and results. Statements occasionally involve some notions that are not introduced formally. In such cases, names are generally chosen to be self-explanatory (for instance, not_wrong_program). In other cases, we forego giving precise definitions as they are not needed to understand the paper (for instance, the types chunk and addressing are unspecified in the definition of state). Our formalization makes an extensive use of inductive definitions, which are introduced in Coq using the keyword Inductive. Inductive definitions are used both for introducing new



Fig. 2: Example program and its SSA forms

datatypes, e.g. the type of RTL instructions in Figure 4, and for introducing inductive relations, e.g. the operational semantics of RTL instructions in Figure 4. In the latter case, the declarations are written according to the pattern

Inductive R : $A \rightarrow B \rightarrow Prop$:= | Rule1: \forall a b, ... \rightarrow R a b | Rule2:... \rightarrow R a b

meaning that the relation R is a binary predicate (indicated by Prop, the type of propositions in Coq) whose arguments are of types A and B respectively. The relation R is defined by two rules Rules1 and Rules2, describing when the proposition (R = b) holds for elements a and b (the hypotheses are indicated by dots).

2. BACKGROUND

2.1. Static Single Assignment form

Static Single Assignment is an intermediate representation in which variables are statically assigned exactly once, thus making explicit in the program syntax the link between the program point where a variable is defined and read.

2.1.1. Converting into SSA form. For straightline code, one simply tags each variable definition with an index, and each variable use with the index corresponding to the last definition of this variable. For example, [x := 1; y := x + 1; x := y - 1; y := x] is transformed into $[x_0 := 1; y_0 := x_0 + 1; x_1 := y_0 - 1; y_1 := x_1]$. The transformation is semantics-preserving, in the sense that the final values of x and y in the first snippet coincide with the final values of x_1 and y_1 in the second snippet.

On the other hand, one cannot transform arbitrary programs into semantically equivalent programs in SSA form solely by tagging variables: one must insert ϕ -functions to handle branching statements. Figure 2 shows a program a), and a program b) that corresponds to a SSA form of a). In program a), the value of variable x read at node 9 either comes from the definition of x at entry or at node 6. In program b), these two definitions of x are renamed into the unique definition of x_0 and x_2 and merged together by the ϕ -function of x_3 at entry of node 9. The precise meaning of a ϕ -block depends on the numbering convention of the predecessor nodes of each junction point. In Figure 2 b) we make explicit this numbering by labelling the CFG edges. For example, node 3 is the first predecessor of point 9 and node 6 is the second one. The semantics of ϕ -functions is given in the seminal paper by Cytron et al. [1991]:

If control reaches node j from its kth predecessor, then the run-time support remembers k while executing the ϕ -functions in j. The value of $\phi(x_1, x_2, ...)$ is just the value of the kth operand. Each execution of a ϕ -function uses only one of the operands, but which one depends on the flow of control just before entering j.

2.1.2. Maximal, minimal and pruned SSA. There may be several SSA forms for a single program CFG. Figure 2 gives alternative SSA forms for a same initial program. In the maximal SSA form (Figure 2b), a ϕ -function is inserted for all program variables, at each join point. As the number of ϕ -functions directly impacts the quality of the subsequent optimizations—as well as the size of the SSA form—it is important that SSA generators for real compilers produce an SSA form with a minimal number of ϕ -functions.

The minimal SSA form is informally specified as follows: a ϕ -function is needed for a given variable at join points that can be reached by at least two distinct definitions of that variable in the initial program. This is captured by the notion of convergence point of CFG paths starting at two distinct definition points of a variable (the join operator in [Cytron et al. 1991]).

Consider the program examples in Figures 2a and 2c. Two definitions of y (at points 1 and 4) can reach the join point 3: a ϕ -instruction is required at node 3 in Program 2c. On the other hand, there is only one definition of x (the initial implicit definition of x) that reaches that point in Program 2a and no ϕ -function is inserted for x at point 3 in Program 2c.

Algorithmically, it is more efficient to determine the placement of ϕ -functions of minimal SSA using the equivalent notion of *dominance frontier*.

Definition 2.1 (Dominance relation). A node i in a CFG dominates another node j if every path from the entry node of the CFG to j contains i. The dominance is said to be strict if additionally $i \neq j$.

Definition 2.2 (Dominance frontier). For a node *i* of a CFG, the dominance frontier DF(i) of *i* is defined as the set of nodes *j* such that *i* dominates at least one predecessor of *j* in the CFG but does not strictly dominate *j* itself. The notion is extended to a set of nodes *S* with $DF(S) = \bigcup_{i \in S} DF(i)$.

Definition 2.3 (Iterated dominance frontier). The iterated dominance frontier $DF^+(S)$ of a set of nodes S is $\lim_{i\to\infty} DF^i(S)$, where $DF^1(S) = DF(S)$ and $DF^{i+1}(S) = DF(S \cup DF^i(S))$.

Efficient algorithms for computing the dominance frontiers rely on an effective representation of the dominance relation, the dominator tree, which in turn relies on the notion of immediate dominator.

Definition 2.4 (Immediate dominator). The immediate dominator of a node j, written idom(j) is the closest strict dominator of j on every path from the entry node to j. It is uniquely determined.

Definition 2.5 (*Dominator tree*). The dominator tree is defined as follows. The start node is the root of the tree. Each node's children are the nodes it immediately dominates.

In a *minimal SSA* program generated by Cytron *et al.*'s algorithm, every ϕ -function of an instance x_i of an original variable x appears in a junction point j if and only if j belongs to the iterated dominance frontier of the set of definition nodes of x in the original program.

However, one can achieve more compact SSA forms by observing that, at any junction point, dead variables need not to be defined by a ϕ -function. The intuition is captured by the notion of *pruned SSA* form: a program is in *pruned SSA* form when the ϕ -functions appear at the iterated dominance frontiers and for each ϕ -function of an instance x_i of an original variable x at a junction point j, x is live at j in the original program (there is a path from j to a use of x that does not redefine x). Compared to

ACM Transactions on Programming Languages and Systems, Vol. Vol, No. Num, Article Art, Publication date: May Year.



Fig. 3: Common sub-expression elimination (CSE) using Gobal Value Numbering (GVN). A no-op instruction is written Inop.

minimal SSA (Figure 2c), pruned SSA detects that the ϕ -function for y at point 9 can be removed. Finally, *semi-pruned* SSA forms can be seen as a good trade-off between a minimal SSA form that is not compact enough, and pruned SSA forms that are sometimes too costly to compute (pathological CFGs with can make the liveness analysis intractable). The liveness analysis used for semi-pruned SSA is local to basic blocks: it is hence less precise, but more efficient.

2.1.3. SSA-based optimizations. The SSA form simplifies the definition of many common optimizations. For instance, copy propagation algorithms can just walk through a SSA program, identify statements of the form x := y, and replace every use of xby y. Furthermore, several optimizations are naturally formulated on SSA. One typical SSA-based optimization is *Global Value Numbering* (GVN) [Alpern et al. 1988], which assigns to variables an identifying number such that variables with the same number will hold equal values at execution time. The effectiveness of GVN lies in its ability to compute efficiently numberings that identify as many variables as possible. Advanced algorithms [Alpern et al. 1988; Briggs et al. 1997] efficiently compute such numberings. We briefly explain one such numbering in Section 7.

Figure 3 illustrates how GVN can be used to eliminate redundant computation. The left program is the original code. In this program, for each *i*, the variables x_i and y_i are assigned the same value number. Hence, the evaluation of $y_1 + 1$ (resp. $y_1 + 2$) is a redundant computation when assigning y_2 (resp. y_3), and one can transform the program into the semantically equivalent one shown on the right of the figure. The strength of the analysis lies in its ability to reason about ϕ -functions, which allows it to infer the equality $x_2 = y_2$. This is only possible because the numbering is global to the whole program. Any block-local analysis would fail to discover the equality $x_2 = y_2$.

2.2. CompCert

CompCert is a realistic, formally verified compiler that generates PowerPC, ARM or x86 code from source programs written in a large subset of C. CompCert formalizes the operational semantics of dozen intermediate languages, and proves a semantics preservation theorem for each phase.

Preservation theorems are expressed in terms of program behaviors, i.e. finite or infinite traces of external function calls (a.k.a. systems calls producing observable events), that are performed during the execution of the program, and claim that individual compilation phases preserve behaviors.

A consequence of the theorems is that for any C program p that does not go wrong (i.e. it does not reach a non-final state where no execution step is valid), and target pro-

gram tp output by the successful compilation of p by the compiler compcert_compiler, the set of behaviors of p contains all behaviors of the target program tp. The formal theorem is:

```
Theorem compcert_compiler_correct: \forall (p: C.program) (tp: Asm.program),
(not_wrong_program p \land compcert_compiler p = OK tp) \rightarrow
(\forall beh, exec_asm_program tp beh \rightarrow exec_C_program p beh).
```

Each phase of the compiler is formally proved relying on simulation techniques, and the formal development of CompCert provides the general correctness theorems of the simulation diagrams. Some parts of the CompCert compiler are not directly proved in Coq. This is the case of the register allocation [Leroy 2009], which is based on a graph coloring algorithm. The graph coloring algorithm is written in OCaml, and then validated a posteriori by a checker written in Coq. The correctness proof of the checker (stating that if a coloring is accepted by the validator, then it is indeed a valid coloring) ensures this compilation phase provides the same guarantees as a transformation written and proved directly in Coq, with the additional benefit of abstracting away complex implementation details and heuristics.

3. THE RTL LANGUAGE

Our middle-end is plugged in at the level of the RTL language in CompCert. This section presents briefly this language. RTL stands for Register Transfer Language. It is a CFG-based three-address-like representation of the code, where most of the existing optimisations are performed (constant propagation, removal of redundant cast, tail call detection, local value numbering and a register allocation that includes copy propagation).

3.1. Syntax and semantics

The syntax and semantics of RTL is given in Figure 4. An RTL program is defined as a set of global variables, a set of functions, and an entry node. Functions are modelled as records that include a function signature fn_sig, and a CFG fn_code of instructions over pseudo-registers. The CFG is not a basic-block graph: instead, it is a partial map from CFG nodes to single instructions (in Figure 4, this map has type PTree.t instr¹), and we stick to this important design choice of CompCert. As explained by Knoop et al. [1998], it allows for simpler implementations of code manipulations and simplifies correctness proofs of analyses or transformations, while still achieving acceptable performance.

The RTL instruction set includes arithmetic operations (Iop), memory loads (Iload) and stores (Istore), function calls (Icall), conditional (Icond) and unconditional jumps (Inop), and a return statement (Ireturn)— we do not discuss here jumptables and other kinds of function calls: calls to a function pointer stored in a register, tail calls, and built-in functions. With the exception of Return, all instructions take a node pc as final argument which denotes the next instruction to execute. Additionally, all instructions but Inop take as arguments pseudo-registers (of type reg), memory chunks, or addressing modes.

The type of states is defined as the tagged union of regular states, call states and return states (Figure 4). We focus on regular states, as we only expose here the intraprocedural part of the language. A regular semantic state (State) is a tuple that contains a call stack (representing the current pending function calls), the current function description and stack pointer (to the stack data block, a part of the global memory

ACM Transactions on Programming Languages and Systems, Vol. Vol, No. Num, Article Art, Publication date: May Year.

Art:8

¹CompCert, and our SSA extension thereof, crucially rely on these data structures. A data structure of type (Ptree.t a) is an associative, partial map where keys have type positive – binary encoding of strictly positive integers – and associated data have type a. As the map is partial, the lookup return type is (option a).

```
Inductive instr :=
                                         RTL instructions (excerpt)
     Inop (pc: node)
     Iop (op: operation) (args: list reg) (res: reg) (pc: node)
    Iload (chk:chunk) (addr:addressing) (args: list reg) (res: reg) (pc: node)
    Istore (chk:chunk) (addr:addressing) (args:list reg) (src: reg) (pc: node)
     Icall (sig: signature) (fn:ident) (args: list reg) (res: reg) (pc: node)
     Icond (cond: condition) (args: list reg) (ifso ifnot: node)
    Ireturn (or: option reg).
Definition code := PTree.t instr. type of code graph
Record function := {
   fn_sig: signature;
                                         function signature
   fn_params: list reg;
                                         parameters
   fn_stacksize: Z;
                                        activation record size
   fn_code: code;
                                         code graph
                                         entry node
   fn_entrypoint: node
}.
Inductive state :=
 | State (stack: list stackframe) call stack
          (f: function)
                                        current function
                                        stack pointer
          (sp: val)
          (pc: node)
                                        current program point
          (rs: regset)
                                        register state
          (m: mem)
                                        memory state
   Callstate (stack: list stackframe) (f: fundef) (args: list val) (m: mem)
   Returnstate (stack: list stackframe) (v: val) (m: mem).
\texttt{Inductive} \texttt{ step: genv} \rightarrow \texttt{state} \rightarrow \texttt{trace} \rightarrow \texttt{state} \rightarrow \texttt{Prop} \texttt{ :=}
 | ex_Inop: \forall ge s f sp pc rs m pc',
       fn_code f pc = Some(Inop pc') \rightarrow
       step ge (State s f sp pc rs m) \epsilon (State s f sp pc' rs m)
 | ex_Iop: \forall ge s f sp pc rs m pc' op args res v,
       fn_code f pc = Some(Iop op args res pc') \rightarrow
       eval_operation sp op (rs##args) m = Some v 
ightarrow
       step ge (State s f sp pc rs m) \epsilon (State s f sp pc' (rs#res\leftarrowv) m)
 | ex_Iload: \forall ge s f sp pc rs m pc' chk addr args res a v,
       fn_code f pc = Some(Iload chk addr args res pc') \rightarrow
       eval_addressing sp addr (rs##args) = Some a \rightarrow Mem.loadv chk m a = Some v \rightarrow
       step ge (State s f sp pc rs m) \epsilon (State s f sp pc' (rs#res\leftarrowv) m)
```

Fig. 4: Syntax and semantics of RTL (excerpt)

where variables dereferenced in the C source program reside), the current program point, the register state (a mapping of local variables to values) and the global memory. The semantics also includes a global environment (of type genv) mapping function names and global variables to memory addresses.

The operational behavior of programs is modelled by the relation step between two semantic states (see Figure 4), and a trace of events. The instructions we show there do not emit any event, hence the transitions that they induce are tagged by the empty event trace ϵ . In the full formalization of RTL, the only instruction producing an event is the call to an external function, which produces for a system call, an event made

ACM Transactions on Programming Languages and Systems, Vol. Vol, No. Num, Article Art, Publication date: May Year.



Fig. 5: An RTL program and its normalized version

of (1) the callee identifier, (2) the values of its arguments, and (3) its return value. We briefly comment on the rules: (Inop pc') branches to the next program point pc'. (Iop op args res pc') performs the arithmetic operation op over the values of registers args (written rs##args), stores the result in res (written $rs#res \leftarrow v$), and branches to pc'. The instruction (Iload chk addr args res pc') loads a chk memory quantity from the address determined by the addressing mode addr and the values of the args registers, stores the memory quantity just read into res, and branches to pc'.

3.2. Normalizing RTL syntax

Before generating the SSA form of an RTL code, we first perform a structural normalization phase of the RTL code (see Figure 1) that we have added to CompCert, prior to the middle-end proper.

This normalization phase consists of transforming an RTL program into another one, with additional structural constraints on the CFGs of functions. We normalize an RTL CFG so that the only instruction that can lead to a junction point is an (Inop pc) instruction. Figure 5 shows an example RTL program and its normalized version.

This normalization phase has been programmed and proved in Coq. One may initially suspect that this normalization phase is quite insignificant. But this structural constraint will carry over the SSA form of RTL programs, and will allow for lightening the formal development of our SSA middle-end. As will be pointed out in the next sections, this impacts the formal definitions of the syntax of SSA, but also greatly simplifies its semantics. This also simplifies the definition and the proof of our SSA validator, the GVN-based CSE, and the SSA deconstruction. Also, we take care during this normalization to remove from the function CFG all the nodes that are not syntactically reachable from the entry node. Having CFGs with all nodes reachable simplifies many formal definitions, including the one of dominance. Strictly speaking, an unreachable node is dominated by any other node in a graph. So, by eliminating such nodes, we also eliminate this corner case from the definitions.

4. THE SSA LANGUAGE

We describe the syntax and operational semantics of our SSA language that provides the SSA form of RTL programs. We equip the notion of SSA program with a *wellformedness* predicate capturing essential properties of SSA forms.

4.1. SSA programs

4.1.1. Syntax. Our definition of SSA program distinguishes between RTL-like instructions and ϕ -functions. The distinction avoids the need for unwieldy mappings between program points when converting to SSA, and allows for a smooth integration in Comp-Cert. Figure 6 introduces the syntax of SSA.

Compared to RTL functions, SSA functions operate on indexed registers of type SSA.reg, and include an additional field fn_phicode mapping junction points to ϕ -

Definition reg := RTL.reg * id	type of indexed registers		
<pre>Inductive instr :=</pre>	RTL-like instructions (operating on SSA.regs)		
<pre>Inductive phiinstr :=</pre>	(res: SSA.reg).	ϕ -functions	
Definition phiblock:= list phi	type of ϕ -blocks		
Definition phicode := PTree.t	phiblock.	type of ϕ -blocks graph	
<pre>Record function := { fn_sig: signature; fn_params: list SSA.reg; fn_stacksize: Z; fn_code: code; fn_phicode: phicode; fn_entrypoint: node }.</pre>	function signature parameters activation record si code graph ϕ -blocks graph entry node	ze	

Fig. 6: Syntax of SSA

blocks. The latter are modelled as lists of ϕ -functions, each of the form (Iphi args res), where res is an indexed register, and args a list of indexed registers.

We define structural constraints that allow giving an intuitive semantics to SSA programs. First, we require that the domain of the function fn_phicode be the set of junction points. Second, we require that all ϕ -functions in a ϕ -block have the same number of arguments as the number of predecessors of that block. Our last requirement is the normalization criterion of the CFG of SSA functions: all predecessors of a junction point must be (Inop pc) instructions.

4.1.2. Strict SSA. We consider two essential properties of SSA forms: unique definitions and strictness [Brisk 2006]. The unique definitions property states that each register is uniquely defined, whereas the strictness property states that each variable use is dominated by the (unique) definition of that variable.

While the two properties are closely related, neither implies the other: the program $[y_0 := x_0; x_0 := 1]$ satisfies the unique definitions property but is not in strict form whereas the program $[x_0 := 1; x_0 := 2; y_0 := x_0]$ is strict but does not satisfy the unique definitions property.

To formalize these properties, one first defines the type of CFG paths, and two predicates dom and sdom for dominance and strict dominance. We also prove many properties of the dominance relation, such as its reflexivity, transitivity, and anti-symmetry. Then, one must define the two predicates def and use of type SSA.function \rightarrow SSA.reg \rightarrow node \rightarrow Prop such that proposition def f x pc (respectively use f x pc) holds iff the register x is defined (resp. used) at node pc in the code of the function f. Predicate def is defined in the obvious way. The definition of use is more involved, because of ϕ -functions. A variable is used either by an RTL-like instruction or a ϕ -function:

Definition use (f:SSA.function) (x:reg) (pc:node) : Prop := use_code f x pc ∨ use_phicode f x pc.

where predicate use_code defines when a variable is used in the RTL-like code. It is defined straightforwardly: a variable is used if it appears in the right hand-side of an assignment, in the condition of an Icond instruction, as an argument of a function



Fig. 7: Example (normalized) SSA program – The variable y_2 is defined at node 3 (in the ϕ -block attached to that program point), and used at points 3 and 4. The variable x_2 is defined at node 6 and used at node 8, the second predecessor of the junction point 9, where x_2 is the second argument of the ϕ -function.

call *etc*. We now explain predicate use_phicode. The widely adopted convention is to view ϕ -functions as lazily evaluated. Hence, the *k*th argument of a ϕ -function is used at the *k*th predecessor of the corresponding block.

```
Inductive use_phicode : SSA.function → reg → node → Prop :=
    | upc_intro : ∀ f pc pred k arg args dst phib
    (PHIB : fn_phicode f pc = Some phib)
    (ASSIG : In (Iphi args dst) phib)
    (KARG : nth_error args k = Some arg) arg is the kth element of args
    (KPRED : index_pred f pred pc = Some k), pred is the kth predecessor of pc in f
    use_phicode f arg pred.
```

This matches the semantics we formally define in Section 4.2: ϕ -functions are executed along the edge leading to the ϕ -block. This definition also allows reusing the traditional notion of strictness defined on non SSA programs. Figure 7 illustrates the definition of predicates def and use.

Using predicates def and use, one can then state the unique definitions and strictness properties, that defines the strict SSA form. We omit the formal definition of unique_def (it is as expected but rather verbose).

```
Definition unique_def (f: SSA.function) := ...

Definition strict (f: SSA.function) :=

\forall x u d, use f x u \rightarrow def f x d \rightarrow dom f d u.
```

4.1.3. Well-formed SSA programs. Finally, the well-formedness of SSA programs is formally defined by the following predicates (the keyword Record must be interpreted as a conjunction):

```
Record wf_ssa_function (f:SSA.function) : Prop := {
fn_ssa: unique_def f;
fn_wf_block: block_nb_args f;
fn_strict: strict f;
fn_block_jp: \forall jp, join_point jp f \leftrightarrow fn_phicode f jp \neq None;
fn_norm:\forall jp pc, join_point jp f \rightarrow jp\in(succs f pc) \rightarrow fn_code f pc=Some(Inop jp)
}.
```

where predicate (join_point jp f) means that the program point jp is a join point in the CFG of function f (meaning it has at least two static predecessors). The predicate

```
Inductive step: SSA.genv \rightarrow SSA.state \rightarrow trace \rightarrow SSA.state \rightarrow Prop :=
   | ex_Inop_njp: ∀ ge s f sp pc rs m pc',
          fn_code f pc = Some(Inop pc') \rightarrow

\neg join_point pc' f \rightarrow
          step ge (State s f sp pc rs m) \epsilon (State s f sp pc' rs m)
    ex_Inop_jp: \forall ge s f sp pc rs m pc' phib k,
         fn_code f pc = Some(nop pc') \rightarrow
fn_phicode f pc' = Some phib\rightarrow
index_pred f pc pc' = Some k \rightarrow
          step ge (State s f sp pc rs m) \epsilon (State s f sp pc' (phistore k rs phib) m)
Fixpoint phistore (k:nat) (rs:SSA.regset) (phib:phiblock) : SSA.regset :=
  match phib with
        \texttt{nil} \Rightarrow \texttt{rs}
        (Iphi args res)::phib \Rightarrow
              match nth_error args k with
                   None \Rightarrow rs
                                                      never happens for well-formed SSA functions
                   Some arg \Rightarrow (phistore k rs phib)#res \leftarrow (rs#arg)
              end
   end.
```

Fig. 8: Semantics of SSA (excerpt)

block_nb_args states that ϕ -functions arguments are consistent with the number of predecessors of the CFG node holding the block. In the sequel, we show that the conversion to SSA yields well-formed programs. Our SSA-based optimizations will assume that the input SSA programs are well-formed. In turn, each of the transformations must be proved to preserve well-formedness. We come back to this point in Section 9.1.

4.2. Semantics

SSA states are similar to RTL states, except that the type of registers and current function are modified into SSA.reg and SSA.function respectively. We describe now the semantics of SSA programs.

4.2.1. Exploiting normalization for an intuitive semantics. The small-step operational semantics is defined on SSA programs that satisfy the structural constraints introduced in the previous section (wf_ssa_function).

Formally, we define SSA.step as a relation between pairs of SSA states and a trace of events. The definition follows the one of RTL.step, except for instructions of the form (Inop pc'), where one distinguishes whether pc' is a junction point or not. In the latter case, the semantics coincide with the RTL semantics, i.e. the program point is updated in the semantic state. If on the contrary pc' is a junction point, then one executes the ϕ -block attached to pc' before the control flows to pc'.

Executing ϕ -blocks on the way to pc' avoids the need to instrument the semantics of SSA with the predecessor program point, and crisply captures the intuitive meaning given to ϕ -blocks by Cytron *et al.* (see Section 2). Note in particular that the normalization ensures that the predecessor of a junction point is an Inop instruction. This greatly simplifies the definition of the semantics (ϕ -block can only be executed after an Inop), and subsequently the proofs about SSA programs.

4.2.2. Parallel execution of ϕ -blocks. Following conventional practice, ϕ -blocks are given a parallel (big-step) semantics. In fact, the SSA generation algorithm ensures, by construction, that the ϕ -functions arguments are never assigned by a distinct ϕ -function in

ACM Transactions on Programming Languages and Systems, Vol. Vol. No. Num, Article Art, Publication date: May Year.

Art:13

the same block. So this parallel semantics seems to be of little help. But later optimizations will exploit this semantics, that makes explicit the independence of ϕ -arguments with regards to ϕ -function destinations [Hack et al. 2006; Boissinot et al. 2009].

The semantics of ϕ -blocks is formally defined with phistore (Figure 8), where we write $S\#r \leftarrow v$ to denote an update of the value of a register r to the value v in a SSA.regset S. When reaching a join point pc' from its kth predecessor, we update the register set rs for each register res assigned in the ϕ -block phib with the value of register arg in rs (written rs#arg), where arg is the kth operand in the ϕ -function of res (written nth_error args k = Some arg). When the ϕ -block is empty (nil clause of the top-level pattern matching), the environment is left unchanged. With the same notations, phistore satisfies, on well-formed SSA functions, a *parallel assignment* property:

```
\forall arg res, In (Iphi args res) phib \rightarrow nth_error args k = Some arg \rightarrow (phistore k rs phib)#res = rs#arg
```

5. TRANSLATION VALIDATION OF SSA GENERATION

Modern compilers typically follow the algorithm by Cytron et al. [1991] to generate a minimal SSA form of programs in almost linear time w.r.t. the size of the program. The algorithm proceeds in four steps:

- (i) Build the CFG dominator tree using the algorithm of [Lengauer and Tarjan 1979]
- (ii) Compute dominance frontiers (bottom-up traversal of the dominator tree)
- (iii) Place ϕ -functions at iterated dominance frontiers of RTL variables
- (iv) Rename definitions and uses of RTL variables with the correct indexes (top-down traversal of the dominator tree).

Programming efficiently the algorithm in Coq and proving formally its correctness is a significant challenge—even verifying formally Step (i) requires one to formalize a substantial amount of graph theory. Instead, we provide a new validation algorithm that checks in linear time that an SSA program is a correct SSA form of an input RTL program. The algorithm is complete w.r.t. minimal SSA form, and can be enhanced by a liveness analysis to handle pruned and semi-pruned SSA forms, as presented in Section 2.1.2. In order to be used in a certified compiler chain, we also show that our validator is sound: it ensure the preservation of behaviors.

Translation validation of SSA conversion is performed in two passes. The first pass performs a structural verification on programs: given a RTL function f and a SSA function tf, it verifies that tf satisfies all clauses of well-formedness except strictness, and that the code of f can be recovered from its SSA form tf simply by erasing ϕ blocks and variable indices—the latter property is captured formally by the proposition structural_spec f tf. The second pass relies on a type system to ensure strictness and semantics-preservation. Overall the pseudo-code of the validator is

where is_well_typed f tf returns true when the function is well-typed with respect to the typing Γ (defined below) in our type system for SSA.

Art:14

ACM Transactions on Programming Languages and Systems, Vol. Vol, No. Num, Article Art, Publication date: May Year.

5.1. Type system

The basic idea of our type system is to track for each variable its *most recent* definition. This is achieved by assigning to all program points a local typing, i.e., an element of $ltype = RTL.reg \rightarrow idx$. We let γ range over local typings. Then, the global typing of an SSA function tf is an element of gtype = node \rightarrow ltype. We let Γ range over global typings. The type system is structured in three layers. The lowest layer checks that RTL-like instructions make a correct use of variables. The middle layer checks that CFG edges are well-typed. Finally, the third layer of the type system defines the notion of well-typed function.

Throughout this section, we use Figure 9 as a running example (an RTL program, its pruned SSA form and its type mapping).

5.1.1. Liveness. As explained in Section 2, liveness information can be used to minimize the number of ϕ -functions in a SSA program. Specifically, ϕ -blocks need to assign only live variables. Hence, our type system is parametrized by a function live modelling a liveness analysis result, a mapping from CFG nodes to sets of registers: (live i) is the set of registers that are live at node i.

Formally, the type system does not need to know much about the liveness information, and how it is computed. We only demand that the live function satisfies two properties: (i) if a variable is used at a program point, then it should be live at this point and (ii) a variable that is live at a given program point is, at the predecessor point, either live or assigned. For a function f, the conjunction of these two properties is denoted by the Coq record (wf_live f live):

```
Record wf_live (f: RTL.function) (live: node \rightarrow Regset.t):= {
   wf_live_use: \forall pc x, use_code f x pc \rightarrow x \in (live pc) ;

   wf_live_incl: \forall pc pc' x,
        is_edge f pc pc' \rightarrow x \in (live pc') \rightarrow
        x \in (live pc) \lor assigned_code f pc x
}.
```

Our type system is able to handle different SSA forms through appropriate instantiations of live. Our formalization provides support for minimal SSA and pruned SSA forms, respectively by defining live as the trivial over-approximation (for each point, it is the set of all the RTL variables), and the result of a standard liveness analysis [Appel 1998a]. One could also support semi-pruned forms, by instantiating live as the result of the block-local liveness analysis of [Briggs et al. 1998]. All these three liveness information can be shown to be well-formed.

Example 5.1 (*Liveness information*). In the last column of the table in Figure 9, we give the liveness information calculated about the variables of the initial RTL function. This information will be used by the validator for validating the pruned SSA form of the program in Figure 9. For instance, the variable y is live at node 3, since it is used at node 3. This variable is however dead (i.e. not live) at point 1 because it is defined (but not used) at this point of the program: it is hence redefined before it is used. At point 6, neither x or y are live. Indeed, the variable x is defined (but not used) at this point 6.

5.1.2. Typing rules for instructions. The type system for instructions checks that RTL-like instructions make a correct use of variables, and that they do not redefine parameters. Its formal definition is given in Figure 10.

Judgments are of the form $\{\gamma\}$ ins $\{\gamma'\}$. Intuitively, the judgment is valid if each variable x is used in ins with the index (γx) , and γ' maps each variable to its last def-



Fig. 9: An RTL program, its pruned SSA form and a valid typing information

inition after execution of ins. The typing rules are formalized as an inductive relation wt_instr. We briefly comment on some rules.

Several rules correspond to instructions that do not define variables, so the input and output local typings are equal. For such rules, one simply checks that the instruction makes a correct use of variables (through use_ok). This is trivially true for instruction (Inop pc) and every local typing γ . The rule for Icond checks that the variables used in the guard are consistent with the local typing input.

In the case of the instruction Iop, which defines the variable (r, i), the output local typing is $\gamma[r \leftarrow i]$, i.e. the input local typing updated for the initial variable r. From this program node onwards, the new version for r is the one indexed with i, and this is the one that should be used later on, until another version for r is defined.

The type system relies on the following convention for function parameters: each of them is given a default index dft (in the example of Figure 9, the default index is 0). The first phase of the validator (check_unique_def) ensures that parameters are not redefined inside the body of the function.

Example 5.2 (*Typing instructions*). We illustrate instruction typing with Figure 9. Consider the input local typing at point 3. The uses of x_0 and y_2 are consistent with it, since $(\Gamma \ 3 \ x) = 0$ and $(\Gamma \ 3 \ y) = 2$. The definition of x_2 at node 6 makes the local typing change for variable x between nodes 6 and 8: it changes from $(\Gamma \ 6 \ x) = 0$ to $(\Gamma \ 8 \ x) = 2$.

5.1.3. Typing rules for edges and functions. The typing rules for edges ensure that ϕ -blocks make a correct use of definitions with regards to a global typing Γ . There are two rules—modelled by the clauses of the inductive relation wt_edge in Figure 10.

The first rule considers the case where the edge does not end in a join point. In this case, typing the edge is equivalent to typing the corresponding instruction.

The second rule considers the case where the edge ends in a junction point: the typing rule checks the ϕ -block attached to it—structural constraints impose that the instruction is an Inop, so we do not need to type-check it. There are three constraints:

- USES ensures that the ϕ -arguments args passed to ϕ -functions are consistent with all incoming local typings: its *k*th argument should be the version of the initial variable brought by the *k*th predecessor of the join point. We omit the formal definition of phiuse_ok
- ASSIG ensures that the output local typing is consistent with the definitions in the ϕ -block
- NASSIG ensures that, if the variable is not assigned in the ϕ -block, then it means that either it is dead, or the incoming indices for this variable are the same for all predecessors

```
Definition use_ok (uses:list SSA.reg)(\gamma:ltype):= \forall r i, In (r,i) uses \rightarrow \gamma r = i.
Inductive wt_instr: ltype \rightarrow SSA.instr \rightarrow ltype \rightarrow Prop :=
| wt_Inop: \forall \gamma s,
        \{\gamma\} Inop s \{\gamma\}
| wt_Istore: \forall~\gamma chk addr args s src,
        use_ok (src::args) \gamma 
ightarrow
        \{\gamma\} Istore chk addr args src s \{\gamma\}
 | wt_Icond: \forall \gamma cond args s1 s2,
        use_ok args \gamma \rightarrow
        \{\gamma\} Icond cond args s1 s2 \{\gamma\}
| wt_Ireturn_some: \forall \ \gamma r,
        \texttt{use\_ok} ~[\texttt{r}] ~\gamma \rightarrow
        \{\gamma\} Ireturn (Some r) \{\gamma\}
| wt_Ireturn_none: \forall \gamma,
        \{\gamma\} Ireturn None \{\gamma\}
| \ \texttt{wt\_Iop:} \ \forall \ \gamma \ \texttt{op args s r i,}
        use_ok args \gamma \rightarrow
        \{\gamma\} Iop op args (r,i) s \{\gamma[r \leftarrow i]\}
 | wt_Iload: \forall~\gamma chk addr args s r i,
        use_ok args \gamma \rightarrow
        \{\gamma\} Iload chk addr args (r,i) s \{\gamma[r \leftarrow i]\}
 | wt_Icall: \forall \ \gamma sig args s id r i,
        use_ok args \gamma \rightarrow
        \{\gamma\} Icall sig id args (r,i) s \{\gamma[r \leftarrow i]\}
Inductive wt_edge (f:SSA.function)(\Gamma:gtype)(live:Regset.t):node \rightarrow node \rightarrow Prop:=
| wt_edge_not_jp: ∀ i j ins
         (NOTJP: fn_code f i = Some ins \land fn_phicode f j = None)
        (WTI: \{\Gamma i\} ins \{\Gamma j\}),
        wt_edge f \Gamma live i j
| wt_edge_jp: ∀ i j ins block
        (JP: fn_code f i = Some ins \land fn_phicode f j = Some block)
        (USES:\forall args r k, In (Iphi args (r,k)) block \rightarrow phiuse_ok r args (preds f j) \Gamma)
        (ASSIG: \forall r k, assigned (r,k) block \rightarrow r \in live \land (\Gamma j r) = k)
        (NASSIG: \forall r, (\forall k, \neg (assigned (r,k) block)) \rightarrow (\Gamma i r = \Gamma j r) \lor r\notinlive),
        wt_edge f \Gamma live i j.
Definition wt_function (f:SSA.function)(\Gamma:gtype)(live:node \rightarrow Regset.t): Prop:=
   \begin{array}{l} (\forall \ i \ j, \ is\_edge \ f \ i \ j \rightarrow wt\_edge \ f \ \Gamma \ (live \ j) \ i \ j) \\ \land \ (\forall \ i \ r, \ fn\_code \ f \ i \ = \ Some \ (Ireturn \ r) \rightarrow \{\Gamma \ i\} \ Ireturn \ r \ \{\Gamma \ i\}) \\ \land \ (\forall \ p, \ In \ p \ (fn\_params \ f) \rightarrow \exists \ r, \ p \ = \ (r, \ \Gamma \ (fn\_entrypoint \ f) \ r) \end{array} 
                                                            \wedge (\Gamma (fn_entrypoint f) r) = dft).
```

Fig. 10: Type system

Art:18

Example 5.3 (*Typing* ϕ -functions). In Figure 9, the ϕ -function for x at point 9 makes correct uses of it because its first argument x_0 matches $(\Gamma 7 x) = 0$ and x_2 matches $(\Gamma 8 x) = 2$. The local typing at node 9 takes into account the definition of x_3 in the block by setting $(\Gamma 9 x)$ to 3. Moreover, no ϕ -function is required for y at node 9 since $y \notin (1ive 9)$, and no ϕ -function is required for x at node 3, since $(\Gamma 2 x) = (\Gamma 5 x)$.

Finally, a function is well-typed with regards to global typing Γ if the local typing induced by Γ at the entry node fn_entrypoint is consistent with the parameters, and all edges and return instructions are well-typed. Return instructions do not correspond to any edge, we thus need to add this constraint explicitly.

5.2. Strictness

All SSA programs accepted by the type system are in strict SSA form. It follows that only well-formed SSA functions will be accepted by the validator.

```
Theorem wt_strict: \forall f tf \Gamma live,
wf_live f live \rightarrow
wt_function f tf \Gamma live \rightarrow
\forall (xi: SSA.reg) (u d: node), use tf xi u \rightarrow def tf xi d \rightarrow dom tf d u.
```

The proof of wt_strict relies on two auxiliary lemmas (explained below) about local typings for well-typed functions. The first lemma states that if a variable x_k is used at node *i*, then it must be that (Γ i x = k).

The second lemma states that, whenever $(\Gamma i x = k)$, the definition point of variable x_k dominates *i*.

Lemma def_gamma : \forall f tf Γ live, wf_live f live \rightarrow wt_function f tf live $\Gamma \rightarrow$ \forall x i d, Γ u x = i \rightarrow def tf (x,i) d \rightarrow dom tf d u.

Under the hypothesis wt_strict, suppose that x_i is used at point u and defined at point d. By use_gamma, we get that $(\Gamma \ u \ x) = i$. We conclude by applying def_gamma to get that d dominates u.

5.3. Soundness

The SSA generation phase, as any other phase of a formally verified compiler must be proved correct in the following sense: all behaviors of the SSA form tf are also behaviors of the corresponding initial RTL function f. In our case, where tf is generated by the untrusted generator and validated a posteriori, we have to prove that if the validator accepts the pair f and tf, then all behaviors of tf are also behaviors of f.

CompCert already provides the general result that a lock-step forward simulation implies preservation of behaviors², it is thus sufficient to exhibit such a simulation, under the assumption that the validator accepts the pair of programs (i.e. all pairs of RTL

 $^{^{2}}$ Technically, a forward simulation implies a backward simulation if the target language is deterministic and the source language is receptive. It is straightforward to prove that SSA is deterministic, but this is not required. Indeed, in CompCert, a forward simulation is first exhibited for each pass of the compilation. Then, they are composed together into a global forward simulation (from deterministic C down to ASM), which implies a backward simulation, thanks to the determinism of the ASM language and the receptiveness of the deterministic variant of C.

and SSA functions in these two programs pass the validator presented in Section 5.1). Such a simulation consists of the three following lemmas:

```
Variable prog:RTL.program.
Variable tprog:SSA.program.
Hypothesis valid_OK : SSA_validator prog tprog = true.
Lemma match_initial_states:
\forall s_1, RTL.initial_state prog s_1 \rightarrow \exists s_2, SSA.initial_state tprog s_2 \land s_1 \simeq s_2.
Lemma match_final_states:
\forall s_1 s_2 r, s_1 \simeq s_2 \rightarrow RTL.final_state s_1 r \rightarrow SSA.final_state s_2 r.
Lemma match_step :
\forall s_1 t s_2, RTL.step (genv prog) s_1 t s_2 \rightarrow \forall s'_1, s_1 \simeq s'_1 \rightarrow \exists s'_2, SSA.step (genv tprog) s'_1 t s'_2 \land s_2 \simeq s'_2.
```

where the binary relation \simeq between semantic states of RTL and SSA carries the invariants needed for proving behavior preservation.

5.3.1. Simulation relation. In particular, \simeq should track the correspondence between the registers of semantics states. To do so, we need to capture the semantics of local typings, that specify the correspondence between the variables of f and tf. This corresponds to the following property:

```
Definition agree (\gamma:ltype) (rs:RTL.regset) (rs':SSA.regset) (live:Regset.t):= \forall r, r \in live \rightarrow rs#r = rs'#(r, \gamma r).
```

This intuitively means that the value of an initial RTL register r is equal to the value of its current version (r, γ) (determined by the local typing γ) in the SSA function. The idea is then to require that, after each computation step, the register states of the RTL and SSA functions agree, with respect to the local typing at the current program point. Note that we will be able to prove such a correspondence only for live variables, and that it is actually sufficient for proving behavior preservation.

Now, defining \simeq only in terms of agreement is not enough to make the proof of simulation go through. We have to constrain more the way RTL and SSA states match. For instance, matching states should have the same memory states and stack pointers. Further, their program counters should be equal. Finally, we add locally to the relation \simeq other invariants relative to the function descriptions of semantic states (e.g. the well-formedness of the SSA function and the well-typedness of the pair of functions).

Formally, the \simeq relation is defined with the inductive match_states below, where we omit, for the sake of brevity, the case for relating semantic states of function calls.

```
Inductive match_states : RTL.state → SSA.state → Prop :=
  | match_states_reg: ∀ s f sp pc rs m ts tf rs' Γ live
    (STACKS: match_stackframes s ts)
    (SPEC: wt_function f tf Γ live)
    (SSA: wf_ssa_function tf)
    (LIVE: wf_live f live)
    (AGREE: agree (Γ pc) rs rs' live),
    (RTL.State s f sp pc rs m) ≃ (SSA.State ts tf sp pc rs' m)
    | match_states_return: ∀ s v m ts
    (STACKS: match_stackframes s ts),
    (RTL.Returnstate s v m) ≃ (SSA.Returnstate ts v m)
    where "s≃t" := (match_states s t).
```

Note that we also define a matching relation for stackframes. This relation basically lifts the invariants of the current functions to the whole callstack of \simeq . This way, at

Art:20

each function call return, the invariants for the caller are available through the matching relation over the stackframes of the callee. This avoids to define (rather clumsily) a global hypothesis on the pair of whole RTL and SSA programs stating the invariants hold for all the functions composing the programs.

5.3.2. Proof sketch. The proof proceeds by nested case-analysis on the kind of semantic state of s1, the relation \simeq , and intruction at the program point under consideration. We treat here the main cases, which are when the instructions are (i) Iop and (ii) Inop when a ϕ -block is attached at its successor point. Consider s1 = (RTL.State sf sp pc rs m) and s1' = (SSA.State ts tf sp pc rs' m), such that (agree (Γ pc) rs rs' (live pc)).

- Suppose (Iop op args res pc') is the instruction at pc in f. Hence, f makes a step towards the state s2 = (RTL.State sf sp pc' (rs#res \leftarrow v) m). By the hypothesis (structural_spec f tf), we know that there is, at point pc in tf, an instruction (Iop op args' (res, i) pc'), and syntax normalization ensures that pc' is not a junction point. Hence, no ϕ -block is attached to it in tf: the matching state is thus s2' = (SSA.State ts tf sp pc' (rs'#(res, i) \leftarrow v) m). In fact both expressions defined by op and respectively args and args' evaluate to the same value v: first, the instruction is well-typed, so that it makes correct uses of its variables, with regards to (Γ pc). Second, rs and rs' agree w.r.t (Γ pc). Finally, all uses are live, by hypothesis on live. Finally, resulting states are still in the relation \simeq , since the update of the local typing specified by the typing rule of the edge (pc, pc') takes into account the actual update of the register states in the semantic step.
- Suppose now (Inop pc') is the instruction at pc in f, with pc' a junction point. In this case, s2 = (RTL.State sf sp pc' rs m). We here take for matching state $s2' = (SSA.State ts tf sp pc' (phi_store kp rs') m)$ where p is the ϕ -block at pc' and k is such that index_pred tf pc pc' = Some k. To show the resulting states stay in the relation, we prove that executing a ϕ -block preserves the agreement between register states (as long at the edge (pc, pc') is well typed. Let x be an RTL variable that is live at pc'. Then, we know that it is live at pc, by the definition of wf_incl and normalization.

If no version of x is assigned in the block, then we use the agreement between rs and rs' at pc. Otherwise, we reason similarly to the case for Iop. We first use hypothesis ASSIG in Figure 10: we have to show that variable x and $(x, (\Gamma pc' x))$ have the same value in the new register states, and this is the case, thanks to constraints we impose on the format of ϕ -blocks, as well as the hypothesis USES in Figure 10: if the kth argument of the ϕ -function is (x, j), then it means that $(\Gamma pc x) = j$, and we can conclude using the agreement of register states at pc.

All other cases are treated similarly in the full formalization, except for executing a function call or return. At function call, we have to prove a partial invariant about the caller (that holds just before calling the function), and the invariants for the callee. The former will then be used at the callee's return, for establishing the rest of the invariant.

5.4. Completeness of the type system

An essential property of our type system is that it accepts all the SSA programs generated by the algorithm by Cytron et al. [1991].

THEOREM 5.4 (TYPE SYSTEM COMPLETENESS). Let f be a normalized RTL function and let tf be the SSA function generated from f by Cytron et al.'s algorithm. Then there exists Γ such that SSA-validator f tf Γ = true.

ACM Transactions on Programming Languages and Systems, Vol. Vol, No. Num, Article Art, Publication date: May Year.

Proving this theorem requires identifying some key properties about the algorithm presented in [Cytron et al. 1991], which we recall in the Section 5.4.1. Given this specification, we show in Section 5.4.2 how to build a global typing, that we prove valid in Section 5.4.3.

This proof is not formalized in the Coq proof assistant. It would require formalizing the specification in Coq, and proving that the actual running algorithm satisfies this specification. Hence, we would not need to run the validator anymore: by the soundness of our type system, we could deduce a full correctness proof of the SSA generation algorithm a la Cytron.

5.4.1. Specification of Cytron et al.'s algorithm. We first review the well-known characterization of the iterated dominance frontier as a fixpoint of the *join* operator J, as well as some properties of the Cytron *et al.*'s algorithm.

Definition 5.5 (Join operator J). Given a set S of nodes, J(S) is defined to be the set of all nodes j such that there are two non empty CFG paths that start at two distinct nodes in S and converge at j, i.e. they both end at j.

LEMMA 5.6 (ITERATED DOMINANCE CHARACTERIZATION). For any set of nodes S, the iterated dominance frontier of S, $DF^+(S)$ satisfies $DF^+(S) = J(S \cup DF^+(S))$.

PROOF. See [Cytron et al. 1991], page 467. □

Let f be an RTL function, and tf the SSA form generated by Cytron's algorithm. For a variable x of f, we write def_x the set of definition points of x in f, and def(x) for the (unique) definition point of the variable in tf. We express now the way Cytron's algorithm defines the set of definition points of the versions of x in tf, and how it determines the right index to use in tf when x is used at some point in f.

LEMMA 5.7 (MINIMAL SSA - DEFINITIONS). Define $D_x = def_x \cup DF^+(def_x)$. D_x is the set of program points where an instance of x is defined in tf, and $DF^+(def_x)$ is the set of nodes where a ϕ -function for x is inserted.

PROOF. Theorem 2 in [Cytron et al. 1991], page 468. □

LEMMA 5.8 (MINIMAL SSA - ABSENCE OF ϕ -FUNCTION). If no instance of a variable x is assigned in the ϕ -block at node n, then a single definition of an instance of x reaches all predecessors of n, without any other instance of x is defined in between.

PROOF. The set of ϕ -functions required for the variable x is by definition $J^+(def_x)$ [Cytron et al. 1991]. We conclude using the definition of the iterated join operator J^+ . \Box

COROLLARY 5.9. If no instance of a variable x is assigned in the ϕ -block at node n, then there exists an instance x_k of x whose definition strictly dominates n.

PROOF. The definition of x_k reaches all predecessors of j, and no instance of x is defined in between. In particular, x_k is defined at a common ancestor of all the predecessors of j. def (x_k) dominates all predecessors of j. It thus dominates j. \Box

LEMMA 5.10 (MINIMAL SSA - USES). If x is used at point i in f, the variable x_k will be used at point i in tf, where x_k is the instance of x such that $def(x_k) \in D_x$ is the closest ancestor of i in the dominator tree of f.

PROOF. See Lemmas 9 and 10 in [Cytron et al. 1991], pages 473-474.

5.4.2. Building a witness global typing. Let f be an RTL function, and tf the SSA form generated by Cytron *et al.*'s algorithm. We explain now how to build a global typing Γ by a depth-first-search (DFS) traversal of the CFG of tf. Each time we reach a new program point j in the DFS, one of its predecessors i in the CFG has already been treated and (Γi) is already defined. To define (Γj), we distinguish two cases:

Case 1 If j is not a join point, for every RTL variable x, we define $(\Gamma j x)$ by case analysis:

— if no instance of x is assigned at *i* in tf, then we set $\Gamma j x = \Gamma i x$;

— if some instance x_k of x is assigned at i in tf, then we set $\Gamma j x = k$;

Case 2 If j is a join point, for every RTL variable x, we define $(\Gamma j x)$ by case analysis on the ϕ -block b at j:

— if no instance of *x* is assigned in *b*, then we set $\Gamma j x = \Gamma i x$;

— if some instance x_k of x is assigned in b then we set $\Gamma j x = k$.

The global typing given in Figure 9 can actually be computed using this construction. Some properties about this witness global typing Γ can be derived, that we will use in the proof of the next paragraph.

LEMMA 5.11 (WITNESS GLOBAL TYPING: PROPERTIES). If $(\Gamma \ i \ x) = k$, then there exists x_k such that $def(x_k)$ dominates i and any shortest CFG path p from $def(x_k)$ to i (excluded) does not go through another definition of an instance of x, i.e. a point in D_x .

PROOF. We proceed by induction on the construction of Γ .

- Base case. For all variables x, (Γ *Entry* x) = dft for all x, and their definition point is the entry point by convention. The condition on shortest paths is trivial since it is empty.
- Induction case. Consider the CFG edge (i, j). We proceed by case analysis on j:
 - . Suppose *j* is not a junction point. By definition of Γ , there are two cases:
 - $(\Gamma j x) = k$ because x_k is defined at *i*. Here, *i* dominates *j*, and the shortest path from *i* to *j* contains only *i* and *j*.
 - $(\Gamma j x) = (\Gamma i x)$ because no instance of x is defined at point i. Applying the induction hypothesis, we get that there is x_k such that $def(x_k)$ dominates i and the shortest path p from $def(x_k)$ to i does not go through another definition of an instance of x. But i dominates j. By transitivity of the dominance relation, we get that $def(x_k)$ dominates j. The minimal path $[def(x_k); \ldots; i; j]$ does not contain any other definition of an instance of x.
 - . Suppose now j is a junction point. There are again two cases.
 - $(\Gamma j x) = k$ because an instance x_k is defined in the ϕ -block at point j. We conclude by the reflexivity of dominance.
 - $(\Gamma \ j \ x) = (\Gamma \ i \ x)$ because no instance of x is defined in the ϕ -block at j. Let $(\Gamma \ i \ x) = k$. Here, the induction hypothesis does not permit to conclude. In this case, j is not in the iterated dominance frontier of any point in def_x (Lemma 5.8). Then, by Corollary 5.9, we get that $def(x_k)$ dominates j.

5.4.3. The witness global typing is a correct typing. Now we prove that tf is typable with Γ as defined in the previous section. We first consider that tf has been generated with a trivial live information full_live, containing at each program point the set of all the RTL variables.

ACM Transactions on Programming Languages and Systems, Vol. Vol, No. Num, Article Art, Publication date: May Year.

Art:22

We consider all edges (i, j) in the CFG of tf, and have to prove that the property $(wf_edge f \Gamma full_live i j)$ holds. We postpone the discussion of typing pruned and semi-pruned SSA versions at the end of the paragraph.

First, we concentrate on verifying that the constraints on the variable definitions are satisfied. We will check that the typing constraints about variables uses (predicate use_ok in Figure 10) in a separated lemma.

LEMMA 5.12 (CONSTRAINTS ON DEFINITIONS). Let (i, j) be an edge in the CFG of tf. Then (wf_edge f Γ full_live i j) holds except for constraints about variable uses.

PROOF. We distinguish two cases.

- Case 1. If j is not a junction point, then i is the sole predecessor of j in the CFG of tf, and (Γj) is defined in terms of (Γi) . In this case, we apply the rule wt_edge_not_jp.
- Case 2. If *j* is a junction point. We have to prove that rule wt_edge_jp is applicable. We consider two cases.
 - Case 2.1. If *i* is the predecessor of *j* in the DFS traversal, Γ *j* is defined in terms of Γ *i*, and the constraints ASSIG and NASSIG hold by definition of Γ. Therefore the edge is typable.
 - Case 2.2. Let i' be the predecessor of j in the DFS, and suppose $i \neq i'$. We have to prove that ASSIG and NASSIG hold. Let b the ϕ -block at point j.
 - ASSIG. Let x_k be assigned in *b*. The live information we use here is full_live, thus *x* is live at point *j*. Additionally, we have $(\Gamma j x) = k$ by construction.

NASSIG. Let x be an RTL variable such that no instance of x is assigned in block b. Because we use full_live, we have to show that $(\Gamma j x) = (\Gamma i x)$.

By definition of Γ , we know that $(\Gamma j x) = (\Gamma i' x)$. It is thus sufficient to prove that $(\Gamma i x) = (\Gamma i' x)$.

If the property would not hold, one could conclude from the Lemma 5.11 that there exist two distinct points ℓ and ℓ' such that a definition of an instance of x occurs in ℓ and ℓ' and there is a path from ℓ (resp. ℓ') that reaches i (resp. i') without meeting any other point in D_x . This implies that $j \in J(D_x) = DF^+(\text{def}_x)$. This leads to a contradiction, as it would mean that j holds a ϕ -node for x (Lemma 5.7). Therefore, an instance of x should be assigned by a ϕ -function in b. This is a contradiction.

This shows that tf is typable with Γ , except for constraints about uses. \Box

LEMMA 5.13 (VARIABLE USES). Let (i, j) be an edge in the CFG of tf. Whenever an instance x_k of x is used at point i in tf, we have $(\Gamma i x = k)$.

PROOF. Suppose that ($\Gamma i x = k'$), with $k' \neq k$. Then, by Lemma 5.11, we know that $def(x_{k'})$ dominates *i*. But x_k is used at point *i*. By Lemma 5.10, we hence know that $def(x_k)$ dominates *i*. Hence, $p_k = def(x_k)$ and $p_{k'} = def(x_{k'})$ both dominate *i*. Therefore, by the property of the dominance relation, either p_k dominates $p_{k'}$ or $p_{k'}$ dominates p_k . We distinguish three cases:

- Case 1. If $p_k = p_{k'}$, we can conclude directly.
- Case 2. Suppose p_k strictly dominates $p_{k'}$. In this case, p'_k would be between p_k and i in the dominator tree. Then, the closest ancestor of i in the dominator tree that belongs to D_x would be $p_{k'}$, and the index used for x at point i should be, by Lemma 5.10, p'_k . This is a contradiction.
- Case 3. Suppose $p_{k'}$ strictly dominates p_k . Then, by antisymmetry of the dominance relation, p_k does not dominate $p_{k'}$. This means that there exists a CFG path p from the entry to $p_{k'}$ that does not go through p_k .

But $(\Gamma i x) = k'$. Thus, by Lemma 5.11, we know it exists a CFG path p' from $p_{k'}$ to i that never meets another point in D_x .

The concatenation of p and p' gives us a path from the entry node of the CFG to i, that never goes through p_k . This contradicts the fact that p_k dominates i.

COROLLARY 5.14 (CONSTRAINTS ON VARIABLE USES). Let (i, j) be an edge in the CFG of tf. Then the constraints on the variable uses in (wf_edge f Γ full_live i j) are satisfied.

Completeness with regards to pruned-SSA form can be shown easily by observing that both the algorithm and the type system make the same use of the liveness information (a dead initial variable does not require a ϕ -function).

5.5. Implementation

For the sake of clarity, we have described a non executable type checker which assumes that structural constraints are satisfied. For efficiency reasons, the Coq implementation of the type system is in fact a bit more complex. In particular, it performs type inference rather than type checking. Additionally, it performs a single, linear scan of the program, and checks the list of arguments of ϕ -functions only once per junction point, rather than once per incoming edge for a given join point.

On the benchmarks given in Section 9.2, our implementation is ten times faster than a type checker derived naively from the non executable type system of Figure 10. We now give an overview of the implementation.

The untrusted SSA generator does not actually compute the whole code of the SSA form of a function. It provides the type checker only with the information that is necessary, namely where to add ϕ -blocks, and how to rename variables definitions in the SSA function – the renaming of variables uses is done directly in Coq. We call this information a hint. It is made of two maps. The first map, of type (PTree.t index), associates to each CFG node (the keys of the map represent positive program points) the index of the variable that this node potentially defines. The second map, of type (PTree.t (PTree.t index)), provides the same kind of information for ϕ -blocks: for a given CFG node (a key in the outter PTree.t), it indicates whenever a block is required (in which case a (PTree.t index) is associated to this key), and in this case, what indexes must be used for the variables defined in that block (here, the keys of the (PTree.t index) denote initial RTL variables, and the associated data are indices). The signature of the external generator for SSA is thus the following:

```
Definition SSA_hint := (PTree.t index) * (PTree.t (PTree.t index)).
Variable extern_SSA_gen: RTL.function \rightarrow (node \rightarrow Regset.t) \rightarrow SSA_hint.
```

Then, given this hint, both the type inference and the code generation (along with the structural checks on the generated code) will be performed simulatenously by the following function:

```
Definition type_infer:
RTL.function \rightarrow (node \rightarrow Regset.t) \rightarrow SSA_hint \rightarrow option SSA.function := ...
```

Since the hint might be incorrect, the type inference may not be able to generate any SSA function, hence the option type of its result. This type inference builds a global typing Γ using the SSA hint, in a way that is similar to the algorithm described in Section 5.4. We then prove that, whenever the inference is successful, the generated function is well typed in the type system described in Figure 10. This is captured by the following theorem:

ACM Transactions on Programming Languages and Systems, Vol. Vol, No. Num, Article Art, Publication date: May Year.

Finally, our SSA generation algorithm is described by the following snippet.

```
Definition ssa_gen (f: RTL.function) : option SSA.function :=
    let live := (LiveAnalysis f) in
    let hint := extern_gen_ssa f live in
    type_infer f live hint.
```

First, a liveness analysis (implemented in Coq) is performed on the RTL function. This liveness information is shared by the external, untrusted SSA generator (written in OCaml) and the type inference. The external SSA generator computes the information (hint) required for the type inference to perform the actual SSA code generation, whilst verifying the validity of the hint.

6. THE SSA EQUATION LEMMA

In this section, we introduce the *equation lemma* that supports the view of programs in SSA form as *systems of equations*. We then illustrate how to reason about a simple SSA-based optimization, namely copy propagation. Using the equation lemma, we will be able in Section 7 to formalize and prove correct a GVN optimization.

6.1. Equation lemma

The SSA representation provides an intuitive reading of programs: one can view the unique definition of a variable as an equation, and by extension one can view SSA programs as systems of equations.

Because every assignment creates a new value name it cannot kill (i.e. invalidate) expressions previously computed from other values. In particular, if two expressions are textually the same, they are sure to evaluate the same result. [Brandis and Mössenböck 1994]

For instance, after at least one iteration, the definitions of x_3 and y_1 respectively induce the two equations $x_3 = y_1 + 1$ (because $x_3 = x_2$ and $x_2 = y_1 + 1$) and $y_1 = x_3 + 1$. There

is however a pitfall: the two equations entail $x_3 = x_3 + 2$, and thus are inconsistent. In fact, equations are only valid at program nodes dominated by the definition that induce them, as captured formally by the *equation-lemma* of SSA:

```
Lemma equation_lemma : ∀ prog d op args x succ f m rs sp pc s,
wf_ssa_program prog →
    reachable prog (State s f sp pc rs m) →
    fn_code f d = Some (Iop op args x succ) →
    sdom f d pc →
    eval_operation sp op (rs##args) m = Some (rs#x).
```

where reachable is a predicate that defines reachable states. In practice, it is often convenient to rely on a corollary that proves the validity of the defining equation of x at program points where x is used – thus avoiding reasoning on the dominance relation. The formal statement of the corollary is obtained by replacing the hypothesis sdom f d pc by the hypothesis use f x pc. The proof of the corollary intensively uses the strictness property of well-formed SSA programs.



6.2. Application with Copy Propagation

We conclude with a succinct account of applying the corollary to prove the soundness of copy propagation (CP)—recall that CP will search for copies x := y and replace every use of x by a use of y. Suppose pc is a program point where such a replacement has been done. Every time pc is reached during the program execution, we are able to derive, using the corollary, that rs#y = rs#x, where rs is the current register state because (i) y is the right hand side of the definition of x and (ii) pc was a use point of x in the initial program. In contrast, on non SSA forms, the reasoning would be more involved since one would have to prove that the reaching definition for x is unique at pc, and that no redefinition of y can occur in between.

7. VALIDATION OF GLOBAL VALUE NUMBERING

This typical SSA-based optimization assigns to variables an identifying number such that variables with the same number will hold equal values at execution time. Several variations of the optimization have been proposed [Alpern et al. 1988; Briggs et al. 1997]. They are generally presented as highly optimised iterative algorithms.

We follow [Alpern et al. 1988] but clearly separate the optimisation into two phases. First, an untrusted analysis, written in OCaml, computes a numbering of SSA programs and for each program point where the numbering detects a redundant computation x := e, it provides a candidate y for replacing the previous operation by x := y. In a second phase, a validator checks the numbering and the proposed assignment simplification.

To achieve this separation of concerns it is useful to reconsider GVN from an abstract interpretation point of view: the analysis computes a fixpoint in the abstract domain of congruence partitions, where partitions are modelled as mappings $\mathcal{N} : \operatorname{reg} \to \operatorname{reg}$ that map a register to the canonical register of its equivalence class (its *number*). The abstract domain is ordered w.r.t. to a partial order \sqsubseteq_{GVN} that coincides with the reverse inclusion of equivalence kernels—recall that the equivalence kernel of \mathcal{N} is the relation $\sim_{\mathcal{N}}$ defined by $\mathbf{x} \sim \mathbf{y}$ if and only if $\mathcal{N} \ x = \mathcal{N} \ y$.

$$\mathcal{N}_1 \sqsubseteq_{\mathrm{GVN}} \mathcal{N}_2 \text{ iff } \sim_{\mathcal{N}_1} \supseteq \sim_{\mathcal{N}_2}$$

The notion of valid numbering is formally defined in Figure 11. First, we define for each numbering \mathcal{N} the relation $\equiv^{\mathcal{N}}$ as the smallest reflexive relation identifying: (i) registers whose assignments share the same operator and corresponding arguments are equivalent w.r.t. \mathcal{N} (predicate same_number) (ii) registers that are defined in the same ϕ -block with equivalent arguments. Then, for a numbering \mathcal{N} to be valid (see GVN_spec), its equivalence kernel must not contain a pair of distinct function parameters and it must moreover be included in $\equiv^{\mathcal{N}}$. The latter ensures the intended post-fixpoint property: if we note n_{param} the numbering that associates each register to itself if it is a function parameter and a default register otherwise, then (GVN_spec \mathcal{N}) is equivalent to $F(\mathcal{N}) \sqsubseteq_{\text{GVN}} \mathcal{N}$ with F the operator defined by $F(\mathcal{N}) = n_{\text{param}} \cap \equiv^{\mathcal{N}}$.

$$\texttt{GVN_spec } \mathcal{N} \text{ iff } n_{\texttt{param}} \cap \equiv^{\mathcal{N}} \ \sqsubseteq_{\texttt{GVN}} \mathcal{N}$$

Viewing the result of the analysis as a post-fixpoint is the key to our second component, a validator that checks whether a numbering \mathcal{N} is indeed a post-fixpoint of the analysis on a program p, and if so returns an optimized SSA program tp. The validator is programmed in Coq, and is accompanied with a proof that optimized programs preserve the behaviors of the original programs.

The crux of the correctness proof of the GVN validator is the correctness lemma for a valid numbering: if \mathcal{N} is a valid numbering for f, and rs is a register state that can be reached at node pc, and x and y are two registers whose definition strictly dominate pc, then $\mathcal{N} = \mathcal{N}$ y entails that rs holds equal values for x and y:

Art:26

ACM Transactions on Programming Languages and Systems, Vol. Vol, No. Num, Article Art, Publication date: May Year.

Fig. 11: Valid numbering

```
Lemma valid_numbering_correct : \forall prog s sp pc rs m,
wf_ssa_program prog \rightarrow GVN_spec \mathcal{N} \rightarrow
reachable prog (State s f sp pc rs m) \rightarrow gamma \mathcal{N} pc rs.
```

where gamma is defined by

```
Definition gamma (\mathcal{N}:reg \rightarrow reg) (pc:node) (rs: regset) : Prop := \forall x y: reg, def_sdom f x pc \rightarrow def_sdom f y pc \rightarrow \mathcal{N} x = \mathcal{N} y \rightarrow rs#x = rs#y.
```

and def_sdom f x pc states that the definition of x in f strictly dominates pc. The definition of def_sdom given below takes care of the case where x is assigned in a ϕ -block at pc (written assigned_phi f pc x). Indeed, a ϕ -block at pc is actually executed before reaching pc while a normal assignment at pc will takes effect after leaving pc.

```
Inductive def_sdom (f:function) (x:reg) (pc:node) : Prop :=
| def_sdom_def_sdom : ∀ def_x,
    def f x def_x → sdom f def_x pc → ¬ assigned_phi f pc x → def_sdom f x pc
| def_sdom_def_phi :
    assigned_phi f pc x → def_sdom f x pc.
```

Let us illustrate the gamma property with Figure 3. Registers x_2 and y_2 share the same numbering: they are indeed equal just after the assignment of y_2 but not before.

Next, we describe the Coq implementation for optimizing SSA programs. The implementation takes as input a numbering \mathcal{N} , and a partial mapping crep that takes as input a register x and node pc and returns, if it exists, a register y such that x and y are related by the equivalence kernel of \mathcal{N} , and the definition of y strictly dominates pc. For efficiency reasons, we do not check the correctness of crep a priori, but lazily during the construction of the optimized program. The optimizer proceeds as follows: first, it checks whether \mathcal{N} satisfies the predicate GVN_spec. Then, for each assignment (Iop op args x pc) of the original SSA program, the optimizer checks whether crep provides a canonical representative y for x at node pc. If so, it checks whether the definition of y strictly dominates pc. This is achieved by means of a dominance analysis, computed directly in Coq with a standard dataflow framework *a la* Kildall³. Provided y is validated, we can safely replace the previous instruction by a move from y to x.

We conclude by commenting briefly on the soundness proof of the transformation. It follows a standard forward simulation proof where the correctness of the numbering

³Computing the dominance analysis with a basic Kildall dataflow can lead to a cubic worst case execution time. We did not encounter any major execution slowdown in our experiments but integrating a more efficient dominance computation would still be valuable for further usage of our SSA middle-end. Recently, Zhao and Zdancewic [2012] have formalized in Coq a faster dominance computation based on Cooper-Harvey-Kennedy algorithm [Cooper et al. 2000]. This non-trivial formalization work could replace advantageously our current approach.

is proved at the same time as the simulation itself. Noticeably, the CFG normalization turned out to be extremely valuable for this proof. Indeed, consider a step from node pc to node pc': we have to prove that (gamma \mathcal{N} pc' rs) holds, asumming (gamma \mathcal{N} pc rs). We reason by case analysis: if the instruction at pc is not an Inop instruction, we know by normalization that pc' is not a junction point. In this case, (def_sdom f x pc') is equivalent to (def_sdom f x pc) \lor (def f x pc) which is particularly useful to exploit the hypothesis that (gamma \mathcal{N} pc rs) holds.

8. CONVERSION OUT OF SSA

The final phase of the middle-end converts SSA programs back to RTL programs, so that they can be further processed by the CompCert back-end, starting with register allocation. Several approaches have been proposed [Sreedhar et al. 1999; Boissinot et al. 2009]. As a first step, we decided to use the conversion described in [Cytron et al. 1991]. The basic idea of this conversion is to substitute each ϕ -function with one variable copy at each predecessor of the junction points:



However, there are several pitfalls to avoid: performing naively the destruction of SSA by such copy insertions can lead to the non preservation of behaviors. Two problems were identified by Briggs et al. [1998]: the presence of critical back-edges (that can lead to the so-called lost-copy problem) and the swap problem. We review both problems in the next sections, and explain how we tackle these two issues. As noted in [Briggs et al. 1998], the swap problem is a particular case of the lost-copy problem but we tackle the issues differently in our development. We finish this section with an overview of the correctness proofs, that shows how the normalization phase can be exploited.

8.1. Critical edges

In the presence of critical back-edges in the progran CFG, the simple copy insertion described above becomes incorrect. We first recall the definition of a critical edge.

Definition 8.1 (Critical edge). A critical edge is an edge (i, j) whose entry *i* has several successors and whose exit *j* has several predecessors.

Figure 12a describes the situation where the exit of the critical edge (i, j) holds a ϕ -block. The problem here is that, the copies cannot be inserted at the predecessors, because they would be executed on some paths that initially did not reach the ϕ -function. This can lead to the well-known lost-copy problem in the presence of critical back-edges and optimizations such as copy folding (see [Briggs et al. 1998]). But copies cannot either be inserted at the edge sink, because it would overwrite the values coming from the others predecessors.

One solution to this problem is to split critical edges, as shown in Figure 12b. After the critical edge (i, j) has been split, the copies for replacing the ϕ -function can be safely inserted at the predecessors of the node holding the block (including the newly inserted node k). Compilers that operate on *basic-block* CFG graphs carefully avoid edge splitting for efficiency concern in later optimization stages. But this is at the cost of making de-SSA algorithms significantly more complex.

Art:28

ACM Transactions on Programming Languages and Systems, Vol. Vol, No. Num, Article Art, Publication date: May Year.



Fig. 12: Critical edges and naive copy insertion

In our case, the normalization we impose on SSA programs pleasingly ensures the absence of critical edges in their CFG. One could fear that the critical edge splitting implied by the normalization could impact later phases of the compiler, but the representation of programs inherited from CompCert deflates this penalty cost. RTL graphs, and thus SSA code graphs, are single-instruction graphs: replacing ϕ -functions with copies automatically splits critical edges by the insertion of code.

8.2. The swap problem

One must also take care of the semantics of ϕ -blocks. They are given a parallel semantics, and, because of optimizations, it is not in general equivalent to a sequential interpretation. Indeed, performing copy propagation on SSA can modify the code, so that ϕ -functions argument and destination registers are no longer independent: a variable x_i can appear both as a source and a target of distinct ϕ -functions in a single ϕ -block. In this case, the copies inserted for converting out of SSA must be sequentialized. This can be done at the reasonable price of inserting at most one temporary variable [Rideau et al. 2008].

In the current state of our development, our conversion out of SSA fails on such ϕ blocks. This is not a limitation in practice, as the GVN optimization we perform on the code does not cause problems of that kind. From the SSA generation until its destruction, the parallel semantics of ϕ -blocks is ensured to be equivalent to the sequential one. We however plan to reuse the work of Rideau et al. [2008] which provides an algorithm for transforming a set of parallel moves into an equivalent sequence of elementary moves (using additional temporaries). This algorithm is already used in CompCert when enforcing calling conventions during the compilation of function calls.

8.3. Correctness proof

For proving the transformation correct, we proceed by giving a forward plus simulation between the SSA program and the RTL program after de-SSA. Indeed, the simulation requires the RTL program to perform one or more steps to simulate the (big-step) execution of a ϕ -block by the initial SSA program.

We also take advantage of the normalization in this proof: the execution of an Inop instruction leading to a junction point with a ϕ -block matches the corresponding inserted copies. Without the normalization, all RTL-like instructions would have resulted in a different case in the proof.

9. IMPLEMENTATION AND EXPERIMENTAL RESULTS

9.1. Coding effort, and ease of proof

We have plugged in CompCert 1.8.2 our SSA middle-end made of (i) a Coq normalization (ii) an OCaml SSA generator and its Coq validator; (iii) an OCaml GVN inference tool and its Coq validator; (iv) a Coq de-SSA transformation. Our formal development

adds 15.000 lines of Coq code and 1.000 lines of OCaml to the 80.000 lines of Coq and 1.000 lines of OCaml provided in CompCert. It does not add any axioms to CompCert.

In terms of development effort, the whole work took approximately 1.5 person year. Special care and successive refinements were necessary to make the verified SSA validator sufficiently efficient, compared to the external, unverified SSA generator. We also spent a significant time in simplifying the formal semantics of the SSA representation, while attempting to build a correctness proof of the GVN optimization. This optimization has hence served has a fruitful *stress test* for our semantic design.

The strong invariants (i.e. strictness, the equation lemma, and the CFG normalization that we propagate from RTL) we have in SSA facilitate optimizations design and proof but, in turn, they must be provably preserved by the transformations. Currently, the SSA optimizations of the middle-end do not modify the CFG of the function, so well-formedness of SSA functions basically boils down to proving that strictness is preserved. In contrast, proving e.g. inlining or aggressive dead-code elimination would be more involved in that respect. For GVN, the proof of strictness preservation is rather trivial. Indeed, when at pc, we replace an operation (Iop op args x pc') by a move instruction (Iop OMove [y] x pc'), we already know that the definition point of the variable y strictly dominates pc. Also, we do not need to prove that the equation lemma is preserved by GVN, as it has been proved to hold for any well-formed SSA function.

In addition to the CFG normalization, the other design choice of our SSA form which has shown to simplify the correctness proof of both the into- and out-of SSA phases is the separate graph of ϕ -blocks. Moreover, this choice appeared to be neutral with regard to the proof of GVN. More generally, we do not expect this ϕ -graph to cause any strong limitation in the implementation (or proof) of other optimizations, as ϕ -blocks must usually be handled as a special case (compared to the rest of the code).

9.2. Experimental evaluation

We use the Coq extraction mechanism to obtain an SSA-based certified compiler, that we evaluate experimentally using the benchmark suite provided with the CompCert distribution. These include around 75.000 lines of C code, and fall into three categories of programs (from 20 to 5.000 LoC): small computation kernels, a raytracer, and the theorem prover Spass⁴. Below we briefly comment on three key points: efficiency of the SSA validator, effectiveness of the GVN optimizer, and efficiency of generated code.

9.2.1. Efficiency of SSA validator. In order to be practical, validators need to be fast enough at compile time. Experimental results are surprisingly good: type-checking the output SSA program is twice as fast as the overall conversion of a program into SSA form. In more detail, the times for SSA generation—specialized to pruned SSA distribute as follows: (i) 9% for normalization of RTL; (ii) 37% for liveness analysis of RTL⁵; (iii) 35% for conversion to SSA using the untrusted OCaml implementation (based on state-of-the-art algorithms); (iv) 19% for validation using the verified validator. This distribution appears to be uniform on all benchmarks except on the biggest functions where the liveness analysis exhibits a non linear complexity.

A possible alternative would be to compute a more efficient (i.e. linear), but less precise, liveness information (such as the one used for building a semi-pruned version of SSA [Briggs et al. 1998]). This could, in theory, lead to a faster computation of liveness, but to an SSA form that is *less* compact, with more variable definitions and uses

 $^{^4}$ Spass is the largest (69.073 LoC), we only use it to evaluate the compilation time.

⁵This analysis, provided in the CompCert distribution, is a traditional backward data-flow analysis based on Kildall's worklist algorithm.

ACM Transactions on Programming Languages and Systems, Vol. Vol, No. Num, Article Art, Publication date: May Year.

	x86				PPC			
	Iop	LVN	GVN	GVN	Iop	LVN	GVN	GVN
				only				only
c. kernels	3,494	163	55	216	3,142	422	54	472
raytracer	2,303	131	29	159	2,755	303	21	322
spass	51,640	122	19	99	$52,\!451$	392	43	306
TOTAL	57,437	416	103	474	58,348	1,117	118	1,100

Table I: **GVN optimizer (x86 and PowerPC backends)** For each set of benchmarks, we count the number of initial Iop in the RTL function (column Iop), the number of Iop optimized away by the LVN-CSE optimization of CompCert (column LVN) and the number of Iop optimized away by our GVN-CSE optimization, right after CompCert's LVN-CSE (column GVN). We also measure the number of Iop that GVN optimizes away without any prior LVN-CSE (column GVN only).

to handle (i.e., in our setting, to both generate and check). During our experiments, pruned SSA appeared to be best trade-off.

Beyond the evaluation of the validator relative to the whole SSA generation phase, it would be interesting to measure the impact of SSA on the whole compilation time. Measuring the into-SSA phase only (in isolation of the out-of-SSA and register allocation phases), would not give much information. However, to date, our de-SSA could be improved: it would not be sufficient if we added new optimizations such as copy propagation, and it introduces too many local registers (which puts the allocator in trouble, as discussed in Section 9.2.3). Improving our de-SSA transformation, or even better, building an allocator on top of SSA – with either techniques similar to [Boissinot et al. 2009] or [Hack et al. 2006] – would allow us to provide relevant measures of the impact of SSA on compilation time. We leave this for future work.

9.2.2. Effectiveness of GVN optimizer. We measure the effectiveness of our GVN analyzer by performing a GVN-based CSE right after a (Local Value Numbering) LVN-based CSE implemented in CompCert. We count how many additional Iop instructions are optimized away by this additional CSE phase. For efficiency concerns about the generated code, we need to keep the LVN phase that eliminates redundant memory loads (currently, this is not done by our GVN optimizer). To keep the comparison fair, we allow CompCert CSE to optimize around function calls—this is disabled in CompCert to keep the register pressure low. The results are given in Table I, for two backends, x86 (left) and PowerPC (right). The overall improvement is significant. Our global CSE optimizes an additional 10% of Iop on PowerPC and an additional 25% on x86.

We also measure how the GVN behaves, without the preliminary LVN optimization. Our global CSE manages to optimize all the Iop instructions that are optimized away by LVN, except 2 for the small computation kernels, and 1 for the raytracer. For Spass, however, GVN only optimizes half the number of Iop. This is due to the fact that in CompCert's LVN, the redundant load elimination and CSE optimizations are interdependent (detecting some redundant loads helps in turn detecting new common sub-expressions, and common sub-expression elimination can lead to extra load redundancy detection).

9.2.3. Efficiency of the Generated code. To assess the efficiency of the generated code, we have compiled the benchmarks with three compilers: CompCert, our version of CompCert extended with a SSA middle-end (CompCertSSA), and gcc - 01. Figure 13 gives the execution times *relative* to CompCert (shorter bars mean faster) on PowerPC.



Fig. 13: Execution times of generated code, normalized to CompCert

As the figure shows, our version of CompCert performs slightly better than CompCert. On average (geometric mean), CompCertSSA code performs almost 2% better than CompCert code (gcc - 01 performs 9% better). The test suite is too small to draw definite conclusions, but the results are encouraging. In particular, adding the SSA middle-end with a simple de-SSA phase does not anihilate the performance of the generated code.

The performance on the integr micro-benchmark seems somewhat anomalous, when comparing with gcc - 01. This can be explained by the fact that CompCert(SSA) does not optimize leaf functions (i.e., functions which make no calls, and hence can run more efficiently if they do not make their own register window), while gcc does. This problem is however somewhat othorgonal to the evaluation of our middle-end.

During our experiments, we observed that the allocator sometimes produces a lot of spill code. The quality of the allocation is impacted mainly by our current SSA deconstruction, that introduces many copies and artificial interferences between variables of a ϕ -block, imposing more constraints on the allocator. Again, we expect that performance would improve significantly by refining our SSA deconstruction, and adapting the register allocator. We leave this interesting challenge for future work. In particular, all the great potential shown by GVN in the above experiment, could be achieved also in terms of execution time.

10. RELATED WORK

We focus on most closely related work and refer to [Leroy 2009] for an overview of mechanized compiler correctness.

10.1. Machine-checked formalizations

Blech et al. [2005] use the Isabelle/HOL proof assistant to verify the generation of machine code from a representation of SSA programs that relies on term graphs. While graph-based representations may be useful for the untrusted parts of our compiler, they increase the complexity of the formal SSA semantics, and make it a greater challenge to verify SSA-based optimizations. They do not provide an algorithm to convert into SSA form, and leave as future work proving the correctness of SSA-based optimizations.

ACM Transactions on Programming Languages and Systems, Vol. Vol, No. Num, Article Art, Publication date: May Year.

A line of work very close to ours is the Vellym project. Zhao et al. [2012] formalize the LLVM SSA intermediate representation in Coq. They define and relate several formal semantics of LLVM, including a static and dynamic semantics. This parallel work shares some similarities with ours concerning the semantics of ϕ -blocks, although they do not normalize the function CFG, nor have a separate graph for ϕ -blocks. Zhao et al. [2013] verify an SSA generation algorithm based on the algorithm of Aycock and Horspool [2000], for the whole LLVM representation. The originality of LLVM's SSA generation is that it starts from a trivial SSA form, where registers are spilled in memory at the end of each basic block. This is intended to simplify the work of LLVM frontends. Then, and this is the hard formalization work, a register promotion algorithm takes care of promoting the so-called *allocas* (i.e. variables local to functions whose addresses can be taken) to registers. The transformation is modelled as a sequence of micro (i.e small and local) transformations, e.g. dead store elimination and dead alloca elimination. Each micro-transformation is then proved correct, and the whole SSA generation is proved correct as a combination of correct micro-transformations. Microtransformations are an interesting and promising approach. In our context, we had more flexibilities. For instance, allocas do not exist in CompCert RTL representation, as they have been eliminated in a prior transformation phase (from Csharpminor to Cminor). Also, we were able to keep a strong correspondence between RTL names and SSA names, and we have kept ϕ -blocks in a distinct code graph, so that the correspondence with the initial RTL program can be easily checked. Another interesting point of their work is the proof methodology they propose for a variety of properties. This is a proof scheme that relaxes the need to consider all variables, thanks to the SSA invariants. Indeed, when reasoning about a program state, only variables definitions that strictly dominate the current program point need to be considered. They apply this proof scheme to prove the equivalent of what we called here the equation lemma (our equation lemma directly embodies the proof scheme), and the strictness property (called the scoping lemma in [Zhao 2013]). In contrast to our work, Zhao et al. [2013]'s SSA generation algorithm is indeed verified in a direct way (and not validated *a posteriori*), but it runs in quadratic time while our generator and its validator run in almost linear time thanks to the Lengauer-Tarjan algorithm. However, as pointed out in [Zhao 2013], their SSA semantics has not yet been used to verify a representative SSA-based optimizations such as GVN, while it constitutes (also in LLVM) a great leverage in generated code performance.

Recently, Schneider [2013] proposed a very elegant approach to formalize SSA in the Coq proof assistant. His work can be construed as a machine-checked formalization of the view that "SSA is functional programming" [Appel 1998b]. Specifically, he considers an intermediate program representation, and equips it with two interpretations, one functional, the other imperative. He analyzes the relationship between both. In particular, he defines a notion of coherent program for which the two semantics coincide, and define back-and-forth translations of programs to their coherent fragment, for both interpretations. In this setting, SSA construction can be understood as turning an imperative program into a coherent, functional one. However, his formalization is only given for a core language, and to our best knowledge, there is not ongoing effort to integrate it into a verified compiler.

Finally, several machine-checked accounts of Continuation Passing Style (CPS) translations exist, e.g. [Dargaye and Leroy 2007; Chlipala 2010], closely related to conversion to SSA form. Chlipala [2010] verifies in Coq a compiler to an idealized as-

sembly language from a small, untyped functional language with mutable references and exceptions. The compiler includes a CPS intermediate representation, on which optimization (CSE) is performed. In contrast to our work, the formalization uses a big-step operational semantics. A great technical achievement of this work is to relax the need of proving the numerous, administrative, and tedious lemmas about substitution. For this, the author relies on the technique of parametric higher-order abstract syntax [Chlipala 2008].

10.2. Translation validation and type systems

Menon et al. [2006] propose a type system that can be used to verify memory safety of programs in SSA form, but their system does not enforce the SSA property. Matsuno and Ohori [2006] define a type system equivalent to SSA: every typable program is given a type annotation making explicit def-use relations. Their type system is similar to ours except they type check one program w.r.t. annotations while we type check a pair of a RTL and a SSA program. They show that common optimizations such as dead code elimination and CSE are type-preserving. But they do not prove the semantics preservation of the optimizations. Stepp et al. [2011] report on a translation validator for LLVM. Their validator uses Equality Saturation [Tate et al. 2009], which views optimizations as equality analyses. Their tool does not validate GVN. Tristan et al. [2011] independently report on an a translation validator for LLVM's inter-procedural optimizations. This tool supports GVN, but is currently not certified.

11. CONCLUSION AND FUTURE WORK

The SSA form is a popular intermediate representation in the compilation community that has been used with great success in many program optimizations since its inception in the late 80's. The structural properties of unique definitions and strictness, as well as the parallel semantics given to ϕ -blocks are the ingredients that led to this success.

If those properties seems rather simple and intuitive, the algorithms underlying the generation of SSA – that actually establish those properties – rely on complex properties of graphs (e.g. the dominator tree or dominance frontiers), that are difficult to justify formally. Moreover, the very semantics of the SSA form has long resisted formalization. As a consequence, the correctness proof of SSA-related algorithms (i.e. generation, optimizations, and destruction), were until very recently not formally proved correct. Over the past few years, some interesting attempts have be made to formalize the semantics of SSA, but these formalizations were rather distant from the intuitive semantics presented in the seminal papers. The correctness of SSA-based analyses and optimizations is usually proved using structural arguments on the CFG only, and the semantic properties and invariants of SSA remain unclear.

In this paper, we have defined a formal semantics for SSA, that is both close to the intuitive definition of the early papers, and amenable for formal reasoning, as witnessed by our fully verified SSA-based middle-end for the verified CompCert C compiler. Thanks to our choices made in the representation of programs, this semantics integrates well in the CompCert architecture. The translation validation approach we use for the conversion to SSA and the GVN optimization allows the middle-end implementing state-of-the-art algorithms, while keeping close to the essence of those phases and to the high-level properties they should satisfy in order to preserve the behaviors of programs. The focused nature of our SSA validator makes it complete with regard to one of the reference implementations of the SSA generators [Cytron et al. 1991], where ϕ -functions placement is determined using dominance-frontiers. We also identified and isolated the semantic counterpart of the structural properties of SSA into a

dedicated invariant lemma (holding at each execution step) on which the correctness proof of several SSA-based optimizations rely.

Our work thus shows that verified compilers can rely on a realistic SSA-based middle-end that implements state-of-the-art algorithms, and opens the way for a new generation of verified compilers based on SSA. Notice that, in general, unverified SSA-based compilers rely on efficient use-def chaining data structures (typically by imperative, destructive constant-time linked-list operations) for implementing optimizations very efficiently. In contrast, we had to implement our middle-end in a pure functional language, in which it is not obvious how to take advantage of SSA's efficient use-def chaining. Implementing such chainings efficiently with functional data-structures is an interesting engineering challenge, although, in our case, no performance bottle-neck was yet encountered on this side.

A priority for further work is to achieve a tighter integration of our middle-end into CompCert. There are three immediate objectives: (i) enhancing our SSA middle-end to handle memory aliases as done by CompCert's RTL-based middle-end, (ii) implementing an SSA-based register allocator [Hack et al. 2006], and (iii) verifying more SSAbased optimizations, including PRE [Chow et al. 1997], or lazy code motion [Knoop et al. 1992]—we expect that our implementation of GVN will provide significant leverage there. Eventually, it should be possible to shift all CompCert optimizations into the SSA middle-end.

ACKNOWLEDGMENTS

We thank Xavier Leroy and the anonymous reviewers for their thoughtful feedback.

REFERENCES

- B. Alpern, M. N. Wegman, and F. K. Zadeck. 1988. Detecting Equality of Variables in Programs. In Proceedings of the 15th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '88). ACM, New York, NY, USA, 1–11. DOI: http://dx.doi.org/10.1145/73560.73561
- A. W. Appel. 1998a. Modern Compiler Implementation: In ML. Cambridge University Press, New York, NY, USA.
- A. W. Appel. 1998b. SSA is Functional Programming. SIGPLAN Notices 33 (1998), 17-20.
- J. Aycock and R. N. Horspool. 2000. Simple Generation of Static Single-Assignment Form. In CC'00 (LNCS), Vol. 1781. Springer Verlag, London, UK, UK, 110–124.
- G. Barthe, D. Demange, and D. Pichardie. 2012. A formally verified SSA-based middle-end Static Single Assignment meets CompCert. In ESOP 2012 (LNCS), Vol. 7211. Springer-Verlag, Berlin, Heidelberg, 47–66.
- J.O. Blech, S. Glesner, J. Leitner, and S. Mülling. 2005. Optimizing Code Generation from SSA Form: A Comparison Between Two Formal Correctness Proofs in Isabelle/HOL. In COCV'05 (ENTCS). Elsevier, Amsterdam, The Netherlands, The Netherlands, 33–51.
- B. Boissinot, A. Darte, F. Rastello, B. Dupont de Dinechin, and C. Guillon. 2009. Revisiting Out-of-SSA Translation for Correctness, Code Quality and Efficiency. In Proc. of the 7th annual IEEE/ACM International Symposium on Code Generation and Optimization (CGO '09). IEEE Computer Society, Washington, DC, USA, 114–125.
- B. Boissinot, S. Hack, D. Grund, B. Dupont de Dinechin, and F. Rastello. 2008. Fast Liveness Checking for SSA form Programs. In Proc. of the 6th annual IEEE / ACM international symposium on Code generation and optimization (CGO '08). ACM, New York, NY, USA, 35–44.
- M. M. Brandis and H. Mössenböck. 1994. Single-pass Generation of Static Single-Assignment Form for Structured Languages. ACM Trans. Program. Lang. Syst. 16, 6 (Nov. 1994), 1684–1698.
- P. Briggs, K.D. Cooper, and L.T. Simpson. 1997. Value Numbering. SPE 27, 6 (1997), 701-724. Issue 6.
- P. Briggs, K. D. Cooper, T. J. Harvey, and L. T. Simpson. 1998. Practical Improvements to the Construction and Destruction of Static Single Assignment Form. *Softw. Pract. Exper.* 28, 8 (July 1998), 859–881.
- P. Brisk. 2006. Advances in Static Single Assignment form and register allocation. Ph.D. Dissertation. University of California at Los Angeles, Los Angeles, CA, USA. Advisor(s) Sarrafzadeh, Majid. AAI3254798.

- A. Chlipala. 2008. Parametric Higher-order Abstract Syntax for Mechanized Semantics. In Proceedings of the 13th ACM SIGPLAN International Conference on Functional Programming (ICFP '08). ACM, New York, NY, USA, 143–156.
- A. Chlipala. 2010. A verified compiler for an impure functional language. In POPL'10. ACM, New York, NY, USA, 93–106.
- F. Chow, S. Chan, R. Kennedy, S-M. Liu, R. Lo, and P. Tu. 1997. A New Algorithm for Partial Redundancy Elimination Based on SSA Form. In Proc. of the ACM SIGPLAN 1997 conference on Programming language design and implementation (PLDI '97). ACM, New York, NY, USA, 273–286.

CompCertSSA 2012. Companion web page. (2012). http://compcertssa.gforge.inria.fr.

- K. D. Cooper, T. J. Harvey, and K. Kennedy. 2000. A simple, fast dominance algorithm. Technical Report. Rice University. Available online at www.cs.rice.edu/~keith/Embed/dom.pdf.
- R. Cytron, J. Ferrante, B.K. Rosen, M.N. Wegman, and F.K. Zadeck. 1991. Efficiently computing static single assignment form and the control dependence graph. ACM TOPLAS 13, 4 (1991), 451–490. Issue 4.
- Z. Dargaye and X. Leroy. 2007. Mechanized Verification of CPS Transformations. In LPAR'07 (LNCS). Springer-Verlag, Berlin, Heidelberg, 211–225.
- S. Hack, D. Grund, and G. Goos. 2006. Register allocation for programs in SSA form. In CC (LNCS). Springer-Verlag, Berlin, Heidelberg, 247–262.
- J. Knoop, D. Koschützkil, and B. Steffen. 1998. Basic-block graphs: Living dinosaurs? In CC. Springer-Verlag, London, UK, UK, 65–79.
- J. Knoop, O. Rüthing, and B. Steffen. 1992. Lazy Code Motion. In PLDI'92. ACM, New York, NY, USA, 224–234.
- T. Lengauer and R.E. Tarjan. 1979. A fast algorithm for finding dominators in a flowgraph. ACM TOPLAS 1, 1 (1979), 121–141. Issue 1.
- X. Leroy. 2009. A Formally Verified Compiler Back-end. JAR 43, 4 (2009), 363-446.
- W. Mansky and E. Gunter. 2010. A Framework for Formal Verification of Compiler Optimizations. In ITP'10. Springer-Verlag, Berlin, Heidelberg, 371–386.
- Y. Matsuno and A. Ohori. 2006. A type system equivalent to static single assignment. In PPDP'06. ACM, New York, NY, USA, 249–260.
- V. Menon, N. Glew, B.R. Murphy, A. McCreight, T. Shpeisman, A.R. Adl-Tabatabai, and L. Petersen. 2006. A verifiable SSA program representation for aggressive compiler optimization. In *POPL'06*. ACM, New York, NY, USA, 397–408.
- G. Necula. 2000. Translation validation for an optimizing compiler. In *PLDI'00*. ACM, New York, NY, USA, 83–94.
- A. Pnueli, M. Siegel, and E. Singerman. 1998. Translation Validation. In TACAS'98 (LNCS). Springer-Verlag, London, UK, UK, 151–166.
- L. Rideau, B.P. Serpette, and X. Leroy. 2008. Tilting at Windmills with Coq: Formal Verification of a Compilation Algorithm for Parallel Moves. JAR 40, 4 (2008), 307–326. Issue 4.
- H. Samet. 1975. Automatically Proving the Correctness of Translations Involving Optimized Code. Ph.D. Dissertation. Computer Science Department, Stanford University. Advisor(s) Cerf, Vinton. STAN-CS-75-498.
- S. Schneider. 2013. Semantics of an Intermediate Language for Program Transformation. Master's thesis. Saarland University.
- V.C. Sreedhar, R. Ju, D.M. Gillies, and V. Santhanam. 1999. Translating Out of Static Single Assignment Form. In SAS'99. Springer-Verlag, London, UK, UK, 194–210.
- M. Stepp, R. Tate, and S. Lerner. 2011. Equality-Based Translation Validator for LLVM. In CAV'11 (LNCS). Springer-Verlag, Berlin, Heidelberg, 737–742.
- R. Tate, M. Stepp, Z. Tatlock, and S. Lerner. 2009. Equality saturation: a new approach to optimization. In *POPL'09*. ACM, New York, NY, USA, 264–276.
- J.B. Tristan, P. Govereau, and G. Morrisett. 2011. Evaluating value-graph translation validation for LLVM. In *PLDI'11*. ACM, New York, NY, USA, 295–305.
- J.B. Tristan and X. Leroy. 2009. Verified validation of lazy code motion. In *PLDI'09*. ACM, New York, NY, USA, 316–326.
- J.B. Tristan and X. Leroy. 2010. A simple, verified validator for software pipelining. In *POPL10*. ACM, New York, NY, USA, 83–92.
- J. Zhao. 2013. Formalizing an SSA-based compiler for verified advanced program transformations. Ph.D. Dissertation. University of Pennsylvania. Advisor(s) Zdancewic, Steve.

ACM Transactions on Programming Languages and Systems, Vol. Vol. No. Num, Article Art, Publication date: May Year.

Art:36

- J. Zhao, S. Nagarakatte, M. Martin, and S. Zdancewic. 2013. Formal verification of SSA-based optimizations for LLVM. In *PLDI'13*. ACM, New York, NY, USA, 175–186.
- J. Zhao and S. Zdancewic. 2012. Mechanized Verification of Computing Dominators for Formalizing Compilers. In CPP'12. Springer, Berlin, Heidelberg, 27–42.
- J. Zhao, S. Zdancewic, S. Nagarakatte, and M. Martin. 2012. Formalizing the LLVM Intermediate Representation for Verified Program Transformation. In *POPL12*. ACM, New York, NY, USA, 427–440.

Received Month Year; revised Month Year; accepted Month Year