

## BIROn - Birkbeck Institutional Research Online

Ströder, T. and Giesl, J. and Brockschmidt, M. and Frohn, F. and Fuhs, Carsten and Hensel, J. and Schneider-Kamp, P. (2014) Proving termination and memory safety for programs with Pointer Arithmetic. In: Demri, S. and Kapur, D. and Weidenbach, C. (eds.) Automated Reasoning. Lecture Notes In Computer Science 8562 8562. New York, U.S.: Springer, pp. 208-223. ISBN 9783319085869.

Downloaded from: <https://eprints.bbk.ac.uk/id/eprint/13532/>

*Usage Guidelines:*

Please refer to usage guidelines at <https://eprints.bbk.ac.uk/policies.html> or alternatively contact [lib-eprints@bbk.ac.uk](mailto:lib-eprints@bbk.ac.uk).

# Proving Termination and Memory Safety for Programs with Pointer Arithmetic<sup>\*</sup>

Thomas Ströder<sup>1</sup>, Jürgen Giesl<sup>1</sup>, Marc Brockschmidt<sup>2</sup>, Florian Frohn<sup>1</sup>,  
Carsten Fuhs<sup>3</sup>, Jera Hensel<sup>1</sup>, and Peter Schneider-Kamp<sup>4</sup>

<sup>1</sup> LuFG Informatik 2, RWTH Aachen University, Germany

<sup>2</sup> Microsoft Research Cambridge, UK

<sup>3</sup> Dept. of Computer Science, University College London, UK

<sup>4</sup> IMADA, University of Southern Denmark, Denmark

**Abstract.** Proving termination automatically for programs with explicit pointer arithmetic is still an open problem. To close this gap, we introduce a novel abstract domain that can track allocated memory in detail. We use it to automatically construct a *symbolic execution graph* that represents all possible runs of the program and that can be used to prove memory safety. This graph is then transformed into an *integer transition system*, whose termination can be proved by standard techniques. We implemented this approach in the automated termination prover AProVE and demonstrate its capability of analyzing C programs with pointer arithmetic that existing tools cannot handle.

## 1 Introduction

Consider the following standard C implementation of `strlen` [23, 30], computing the length of the string at pointer `str`. In C, strings are usually represented as a pointer `str` to the heap, where all following memory cells up to the first one that contains the value 0 are allocated memory and form the value of the string.

```
int strlen(char* str) {char* s = str; while(*s) s++; return s-str;}
```

To analyze algorithms on such data, one has to handle the interplay between addresses and the values they point to. In C, a violation of *memory safety* (e.g., dereferencing `NULL`, accessing an array outside its bounds, etc.) leads to undefined behavior, which may also include non-termination. Thus, to prove termination of C programs with low-level memory access, one must also ensure memory safety. The `strlen` algorithm is memory safe and terminates because there is some address `end`  $\geq$  `str` (an *integer property* of `end` and `str`) such that `*end` is 0 (a *pointer property* of `end`) and all addresses `str`  $\leq$  `s`  $\leq$  `end` are allocated. Other typical programs with pointer arithmetic operate on arrays (which are just sequences of memory cells in C). In this paper, we present a novel approach to prove memory safety and termination of algorithms on integers and pointers automatically. To avoid handling the intricacies of C, we analyze programs in the platform-independent intermediate representation (IR) of the LLVM compilation framework [17]. Our approach works in three steps: First, a *symbolic execution graph* is created

---

<sup>\*</sup> Supported by DFG grant GI 274/6-1 and Research Training Group 1298 (*AlgoSyn*).

that represents an over-approximation of all possible program runs. We present our abstract domain based on *separation logic* [22] and the automated construction of such graphs in Sect. 2. In this step, we handle all issues related to memory, and in particular prove memory safety of our input program. In Sect. 3, we describe the second step of our approach, in which we generate an *integer transition system* (ITS) from the symbolic execution graph, encoding the essential information needed to show termination. In the last step, existing techniques for integer programs are used to prove termination of the resulting ITS. In Sect. 4, we compare our approach with related work and show that our implementation in the termination prover AProVE proves memory safety and termination of typical pointer algorithms that could not be handled by other tools before.

## 2 From LLVM to Symbolic Execution Graphs

In Sect. 2.1, we introduce concrete LLVM states and *abstract* states that represent *sets* of concrete states, cf. [9]. Based on this, Sect. 2.2 shows how to construct symbolic execution graphs automatically. Sect. 2.3 presents our algorithm to *generalize* states, needed to always obtain *finite* symbolic execution graphs.

To simplify the presentation, we restrict ourselves to a single LLVM function without function calls and to types of the form *in* (for *n*-bit integers), *in\** (for pointers to values of type *in*), *in\*\**, *in\*\*\**, etc. Like many other approaches to termination analysis, we disregard integer overflows and assume that variables are only instantiated with signed integers appropriate for their type. Moreover, we assume a 1 byte data alignment (i.e., values may be stored at any address).

### 2.1 Abstract Domain

Consider the `strlen` function from Sect. 1. In the corresponding LLVM code,<sup>5</sup> `str` has the type `i8*`, since it is a pointer to the string’s first character (of type `i8`). The program is split into the *basic blocks* `entry`, `loop`, and `done`. We will explain this LLVM code in detail when constructing the symbolic execution graph in Sect. 2.2.

```
define i32 @strlen(i8* str) {
entry: 0: c0 = load i8* str
      1: c0zero = icmp eq i8 c0, 0
      2: br i1 c0zero, label done, label loop
loop:  0: olds = phi i8* [str,entry],[s,loop]
      1: s = getelementptr i8* olds, i32 1
      2: c = load i8* s
      3: czero = icmp eq i8 c, 0
      4: br i1 czero, label done, label loop
done:  0: sfin = phi i8* [str,entry],[s,loop]
      1: sfinint = ptrtoint i8* sfin to i32
      2: strint = ptrtoint i8* str to i32
      3: size = sub i32 sfinint, strint
      4: ret i32 size }
```

Concrete LLVM states consist of the program counter, the values of local variables, and the state of the memory. The program counter is a 3-tuple  $(\mathbf{b}_{prev}, \mathbf{b}, i)$ , where  $\mathbf{b}$  is the name of the current basic block,  $\mathbf{b}_{prev}$  is the previously executed

<sup>5</sup> This LLVM program corresponds to the code obtained from `strlen` with the Clang compiler [8]. To ease readability, we wrote variables without “%” in front (i.e., we wrote “`str`” instead of “`%str`” as in proper LLVM) and added line numbers.

block,<sup>6</sup> and  $i$  is the index of the next instruction. So if  $Blks$  is the set of all basic blocks, then the set of code positions is  $Pos = (Blks \cup \{\varepsilon\}) \times Blks \times \mathbb{N}$ . We represent assignments to the local program variables  $\mathcal{V}_{\mathcal{P}}$  (e.g.,  $\mathcal{V}_{\mathcal{P}} = \{\mathbf{str}, \mathbf{c0}, \dots\}$ ) as functions  $s : \mathcal{V}_{\mathcal{P}} \rightarrow \mathbb{Z}$ . The state of the memory is represented by a partial function  $m : \mathbb{N}_{>0} \rightarrow \mathbb{Z}$  with finite domain that maps addresses to integer values. So a concrete LLVM state is a 3-tuple  $(p, s, m) \in Pos \times (\mathcal{V}_{\mathcal{P}} \rightarrow \mathbb{Z}) \times (\mathbb{N}_{>0} \rightarrow \mathbb{Z})$ .

To model violations of memory safety, we introduce a special state  $ERR$  to be reached when accessing non-allocated memory. So  $(p, s, m)$  denotes only memory safe states where all addresses in  $m$ 's domain are allocated. Let  $\rightarrow_{\text{LLVM}}$  be LLVM's evaluation relation on concrete states, i.e.,  $(p, s, m) \rightarrow_{\text{LLVM}} (\bar{p}, \bar{s}, \bar{m})$  holds iff  $(p, s, m)$  evaluates to  $(\bar{p}, \bar{s}, \bar{m})$  by executing one LLVM instruction. Similarly,  $(p, s, m) \rightarrow_{\text{LLVM}} ERR$  means that the instruction at position  $p$  accesses an address where  $m$  is undefined. An LLVM program is *memory safe* for  $(p, s, m)$  iff there is no evaluation  $(p, s, m) \rightarrow_{\text{LLVM}}^+ ERR$ , where  $\rightarrow_{\text{LLVM}}^+$  is the transitive closure of  $\rightarrow_{\text{LLVM}}$ .

To formalize *abstract* states that stand for sets of concrete states, we use a fragment of *separation logic* [22]. Here, an infinite set of symbolic variables  $\mathcal{V}_{\text{sym}}$  with  $\mathcal{V}_{\text{sym}} \cap \mathcal{V}_{\mathcal{P}} = \emptyset$  can be used in place of concrete integers. We represent abstract states as tuples  $(p, LV, KB, AL, PT)$ . Again,  $p \in Pos$  is the program counter. The function  $LV : \mathcal{V}_{\mathcal{P}} \rightarrow \mathcal{V}_{\text{sym}}$  maps every local variable to a symbolic variable. To ease the generalization of states in Sect. 2.3, we require injectivity of  $LV$ . The *knowledge base*  $KB \subseteq QF\_IA(\mathcal{V}_{\text{sym}})$  is a set of pure quantifier-free first-order formulas that express integer arithmetic properties of  $\mathcal{V}_{\text{sym}}$ .

The *allocation list*  $AL$  contains expressions of the form  $alloc(v_1, v_2)$  for  $v_1, v_2 \in \mathcal{V}_{\text{sym}}$ , which indicate that  $v_1 \leq v_2$  and that all addresses between  $v_1$  and  $v_2$  are allocated. Finally,  $PT$  is a set of “points-to” atoms  $v_1 \hookrightarrow_{\mathbf{ty}} v_2$  where  $v_1, v_2 \in \mathcal{V}_{\text{sym}}$  and  $\mathbf{ty}$  is an LLVM type. This means that the value  $v_2$  of type  $\mathbf{ty}$  is stored at the address  $v_1$ . Let  $size(\mathbf{ty})$  be the number of bytes required for values of type  $\mathbf{ty}$  (e.g.,  $size(\mathbf{i8}) = 1$  and  $size(\mathbf{i32}) = 4$ ). As each memory cell stores one byte,  $v_1 \hookrightarrow_{\mathbf{i32}} v_2$  means that  $v_2$  is stored in the four cells at the addresses  $v_1, \dots, v_1 + 3$ .

**Definition 1 (Abstract States).** Abstract states *have the form*  $(p, LV, KB, AL, PT)$  where  $p \in Pos$ ,  $LV : \mathcal{V}_{\mathcal{P}} \rightarrow \mathcal{V}_{\text{sym}}$  is injective,  $KB \subseteq QF\_IA(\mathcal{V}_{\text{sym}})$ ,  $AL \subseteq \{alloc(v_1, v_2) \mid v_1, v_2 \in \mathcal{V}_{\text{sym}}\}$ , and  $PT \subseteq \{(v_1 \hookrightarrow_{\mathbf{ty}} v_2) \mid v_1, v_2 \in \mathcal{V}_{\text{sym}}, \mathbf{ty} \text{ is an LLVM type}\}$ . Additionally, there is a state  $ERR$  for violations of memory safety.

We often identify  $LV$  with the set of equations  $\{\mathbf{x} = LV(\mathbf{x}) \mid \mathbf{x} \in \mathcal{V}_{\mathcal{P}}\}$  and extend  $LV$  to a function from  $\mathcal{V}_{\mathcal{P}} \uplus \mathbb{Z}$  to  $\mathcal{V}_{\text{sym}} \uplus \mathbb{Z}$  by defining  $LV(z) = z$  for all  $z \in \mathbb{Z}$ . As an example, consider the following abstract state for our `strlen` program:

$$(\varepsilon, \mathbf{entry}, 0), \left( \begin{array}{l} \{\mathbf{str} = u_{\mathbf{str}}, \dots, \mathbf{size} = u_{\mathbf{size}}\}, \quad \{z = 0\}, \\ \{alloc(u_{\mathbf{str}}, v_{\mathbf{end}})\}, \quad \{v_{\mathbf{end}} \hookrightarrow_{\mathbf{i8}} z\} \end{array} \right). \quad (\dagger)$$

It represents states at the beginning of the `entry` block, where  $LV(\mathbf{x}) = u_{\mathbf{x}}$  for all  $\mathbf{x} \in \mathcal{V}_{\mathcal{P}}$ , the memory cells between  $LV(\mathbf{str}) = u_{\mathbf{str}}$  and  $v_{\mathbf{end}}$  are allocated, and the value at the address  $v_{\mathbf{end}}$  is  $z$  (where the knowledge base implies  $z = 0$ ).

To define the semantics of abstract states  $a$ , we introduce the formulas  $\langle a \rangle_{SL}$  and  $\langle a \rangle_{FO}$ . The separation logic formula  $\langle a \rangle_{SL}$  defines which concrete states are

<sup>6</sup>  $\mathbf{b}_{prev}$  is needed for `phi` instructions (cf. Sect. 2.2). In the beginning, we set  $\mathbf{b}_{prev} = \varepsilon$ .

represented by  $a$ . The first-order formula  $\langle a \rangle_{FO}$  is used to construct symbolic execution graphs, allowing us to use standard SMT solving for all reasoning in our approach. Moreover, we also use  $\langle a \rangle_{FO}$  for the subsequent generation of integer transition systems from the symbolic execution graphs. In addition to  $KB$ ,  $\langle a \rangle_{FO}$  states that the expressions  $alloc(v_1, v_2) \in AL$  represent disjoint intervals and that two addresses must be different if they point to different values in  $PT$ .

In  $\langle a \rangle_{SL}$ , we combine the elements of  $AL$  with the separating conjunction “ $*$ ” to ensure that different allocated memory blocks are disjoint. Here, as usual  $\varphi_1 * \varphi_2$  means that  $\varphi_1$  and  $\varphi_2$  hold for disjoint parts of the memory. In contrast, the elements of  $PT$  are combined by the ordinary conjunction “ $\wedge$ ”. So  $v_1 \hookrightarrow_{\mathbf{ty}} v_2 \in PT$  does not imply that  $v_1$  is different from other addresses occurring in  $PT$ . Similarly, we also combine the two formulas resulting from  $AL$  and  $PT$  by “ $\wedge$ ”, as both express different properties of memory addresses.

**Definition 2 (Representing States by Formulas).** For  $v_1, v_2 \in \mathcal{V}_{sym}$ , let  $\langle alloc(v_1, v_2) \rangle_{SL} = v_1 \leq v_2 \wedge (\forall x. \exists y. (v_1 \leq x \leq v_2) \Rightarrow (x \hookrightarrow y))$ . Due to the two’s complement representation, for any LLVM type  $\mathbf{ty}$ , we define  $\langle v_1 \hookrightarrow_{\mathbf{ty}} v_2 \rangle_{SL} =$

$$\langle v_1 \hookrightarrow_{size(\mathbf{ty})} v_3 \rangle_{SL} \wedge (v_2 \geq 0 \Rightarrow v_3 = v_2) \wedge (v_2 < 0 \Rightarrow v_3 = v_2 + 2^{8 \cdot size(\mathbf{ty})}),$$

where  $v_3 \in \mathcal{V}_{sym}$  is fresh. Here,<sup>7</sup>  $\langle v_1 \hookrightarrow_0 v_3 \rangle_{SL} = true$  and  $\langle v_1 \hookrightarrow_{n+1} v_3 \rangle_{SL} = v_1 \hookrightarrow (v_3 \bmod 256) \wedge \langle (v_1 + 1) \hookrightarrow_n (v_3 \text{ div } 256) \rangle_{SL}$ . Then  $a = (p, LV, KB, AL, PT)$  is represented by<sup>8</sup>  $\langle a \rangle_{SL} = LV \wedge KB \wedge (*_{\varphi \in AL} \langle \varphi \rangle_{SL}) \wedge (\bigwedge_{\varphi \in PT} \langle \varphi \rangle_{SL})$ .

Moreover, the following first-order information on  $\mathcal{V}_{sym}$  is deduced from an abstract state  $a = (p, LV, KB, AL, PT)$ . Let  $\langle a \rangle_{FO}$  be the smallest set with

$$\begin{aligned} \langle a \rangle_{FO} = & KB \cup \{v_1 \leq v_2 \mid alloc(v_1, v_2) \in AL\} \cup \\ & \{v_2 < w_1 \vee w_2 < v_1 \mid alloc(v_1, v_2), alloc(w_1, w_2) \in AL, (v_1, v_2) \neq (w_1, w_2)\} \cup \\ & \{v_1 \neq w_1 \mid (v_1 \hookrightarrow_{\mathbf{ty}} v_2), (w_1 \hookrightarrow_{\mathbf{ty}} w_2) \in PT \text{ and } \models \langle a \rangle_{FO} \Rightarrow v_2 \neq w_2\}. \end{aligned}$$

Let  $\mathcal{T}(\mathcal{V}_{sym})$  be the set of all arithmetic terms containing only variables from  $\mathcal{V}_{sym}$ . Any function  $\sigma : \mathcal{V}_{sym} \rightarrow \mathcal{T}(\mathcal{V}_{sym})$  is called an *instantiation*. Thus,  $\sigma$  does not instantiate  $\mathcal{V}_{\mathcal{P}}$ . Instantiations are extended to formulas in the usual way, i.e.,  $\sigma(\varphi)$  instantiates every  $v \in \mathcal{V}_{sym}$  that occurs free in  $\varphi$  by  $\sigma(v)$ . An instantiation is called *concrete* iff  $\sigma(v) \in \mathbb{Z}$  for all  $v \in \mathcal{V}_{sym}$ . Then an abstract state  $a$  at position  $p$  represents those concrete states  $(p, s, m)$  where  $(s, m)$  is a *model* of  $\sigma(\langle a \rangle_{SL})$  for a concrete instantiation  $\sigma$  of the symbolic variables. So for example, the abstract state  $(\dagger)$  on the previous page represents all concrete states  $((\varepsilon, \mathbf{entry}, 0), s, m)$  where  $m$  is a memory that stores a string at the address  $s(\mathbf{str})$ .<sup>9</sup>

<sup>7</sup> We assume a little-endian data layout (where least significant bytes are stored in the lowest address). A corresponding representation could also be defined for big-endian layout. This layout information is necessary to decide which concrete states are represented by abstract states, but it is not used when constructing symbolic execution graphs (i.e., our remaining approach is independent of such layout information).

<sup>8</sup> We identify *sets* of first-order formulas  $\{\varphi_1, \dots, \varphi_n\}$  with their conjunction  $\varphi_1 \wedge \dots \wedge \varphi_n$ .

<sup>9</sup> The reason is that then there is an address  $end \geq s(\mathbf{str})$  such that  $m(end) = 0$  and  $m$  is defined for all numbers between  $s(\mathbf{str})$  and  $end$ . Hence,  $(s, m) \models \sigma(\langle a \rangle_{SL})$  holds for an instantiation with  $\sigma(u_x) = s(\mathbf{x})$  for all  $\mathbf{x} \in \mathcal{V}_{\mathcal{P}}$ ,  $\sigma(v_{end}) = end$ , and  $\sigma(z) = 0$ .

It remains to define when  $(s, m)$  is a *model* of a formula from our fragment of separation logic. For  $s : \mathcal{V}_{\mathcal{P}} \rightarrow \mathbb{Z}$  and any formula  $\varphi$ , let  $s(\varphi)$  result from replacing all  $\mathbf{x} \in \mathcal{V}_{\mathcal{P}}$  in  $\varphi$  by  $s(\mathbf{x})$ . Note that by construction, local variables  $\mathbf{x}$  are never quantified in our formulas. Then we define  $(s, m) \models \varphi$  iff  $m \models s(\varphi)$ .

We now define  $m \models \psi$  for formulas  $\psi$  that may still contain symbolic variables from  $\mathcal{V}_{sym}$  (this is needed for Sect. 2.2). As usual, all free variables  $v_1, \dots, v_n$  in  $\psi$  are implicitly universally quantified, i.e.,  $m \models \psi$  iff  $m \models \forall v_1, \dots, v_n. \psi$ . The semantics of arithmetic operations and relations and of first-order connectives and quantifiers is as usual. In particular, we define  $m \models \forall v. \psi$  iff  $m \models \sigma(\psi)$  holds for all instantiations  $\sigma$  where  $\sigma(v) \in \mathbb{Z}$  and  $\sigma(w) = w$  for all  $w \in \mathcal{V}_{sym} \setminus \{v\}$ .

We still have to define the semantics of  $\hookrightarrow$  and  $*$  for variable-free formulas. For  $z_1, z_2 \in \mathbb{Z}$ , let  $m \models z_1 \hookrightarrow z_2$  hold iff  $m(z_1) = z_2$ .<sup>10</sup> The semantics of  $*$  is defined as usual in separation logic: For two partial functions  $m_1, m_2 : \mathbb{N}_{>0} \rightarrow \mathbb{Z}$ , we write  $m_1 \perp m_2$  to indicate that the domains of  $m_1$  and  $m_2$  are disjoint and  $m_1 \cdot m_2$  denotes the union of  $m_1$  and  $m_2$ . Then  $m \models \varphi_1 * \varphi_2$  iff there exist  $m_1 \perp m_2$  such that  $m = m_1 \cdot m_2$  where  $m_1 \models \varphi_1$  and  $m_2 \models \varphi_2$ .

As usual, “ $\models \varphi$ ” means that  $\varphi$  is a tautology, i.e., that  $(s, m) \models \varphi$  holds for any  $s : \mathcal{V}_{\mathcal{P}} \rightarrow \mathbb{Z}$  and  $m : \mathbb{N}_{>0} \rightarrow \mathbb{Z}$ . Clearly,  $\models \langle a \rangle_{SL} \Rightarrow \langle a \rangle_{FO}$ , i.e.,  $\langle a \rangle_{FO}$  contains first-order information that holds in every concrete state represented by  $a$ .

## 2.2 Constructing Symbolic Execution Graphs

We now show how to automatically generate a *symbolic execution graph* that over-approximates all possible executions of a given program. For this, we present symbolic execution rules for some of the most important LLVM instructions. Other instructions can be handled in a similar way, cf. [26]. Note that in contrast to other formalizations of LLVM’s operational semantics [31], our rules operate on *abstract* instead of concrete states to allow a *symbolic* execution of LLVM. In particular, we also have rules for refining and generalizing abstract states.

Our analysis starts with the set of initial states that one wants to analyze for termination, e.g., all states where `str` points to a *string*. So in our example, we start with the abstract state  $(\dagger)$ . Fig. 1 depicts the symbolic execution graph for `strlen`. Here, we omitted the component  $AL = \{alloc(u_{str}, v_{end})\}$ , which stays the same in all states in this example. We also abbreviated parts of  $LV, KB, PT$  by “...”. Instead of  $v_{end} \hookrightarrow_{i8} z$  and  $z = 0$ , we directly wrote  $v_{end} \hookrightarrow 0$ , etc.

The function `strlen` starts with loading the character at address `str` to `c0`. Let  $p:ins$  denote that *ins* is the instruction at position  $p$ . Our first rule handles the case  $p$ : “`x = load ty* ad`”, i.e., the value of type `ty` at the address `ad` is assigned to the variable `x`. In our rules, let  $a$  always denote the abstract state *before* the execution step (i.e., above the horizontal line of the rule). Moreover, we write  $\langle a \rangle$  instead of  $\langle a \rangle_{FO}$ . As each memory cell stores one byte, in the `load`-rule we first have to check whether the addresses `ad, \dots, ad + size(ty) - 1` are allocated, i.e., if there is an  $alloc(v_1, v_2) \in AL$  such that  $\langle a \rangle \Rightarrow (v_1 \leq LV(\mathbf{ad}) \wedge$

<sup>10</sup> We use “ $\hookrightarrow$ ” instead of “ $\mapsto$ ” in separation logic, since  $m \models z_1 \mapsto z_2$  would imply that  $m(z)$  is undefined for all  $z \neq z_1$ . This would be inconvenient in our formalization, since  $PT$  usually only contains information about a *part* of the allocated memory.

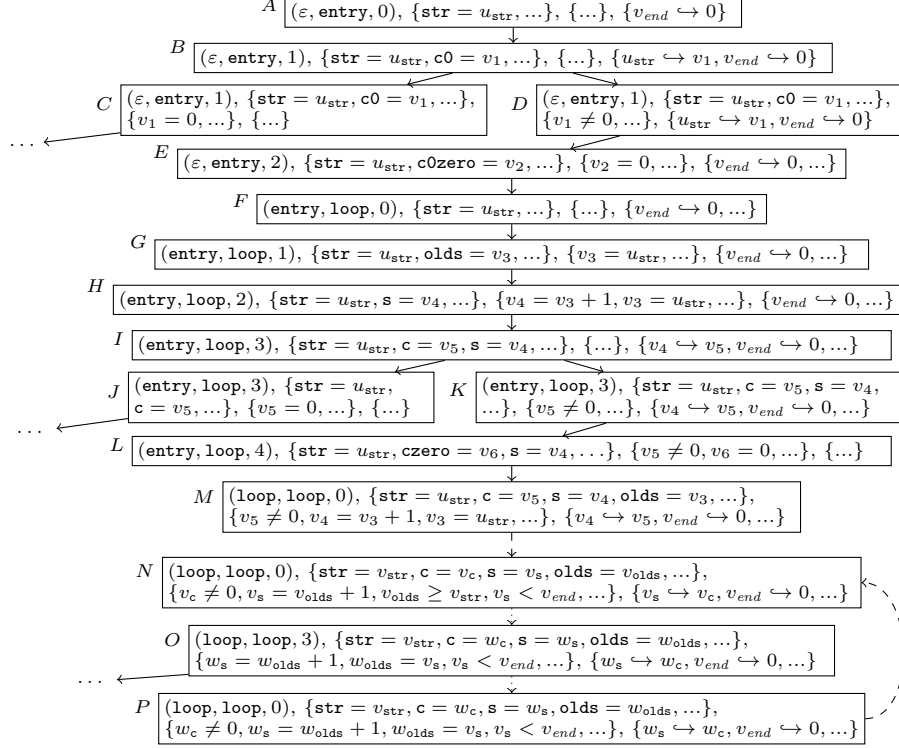
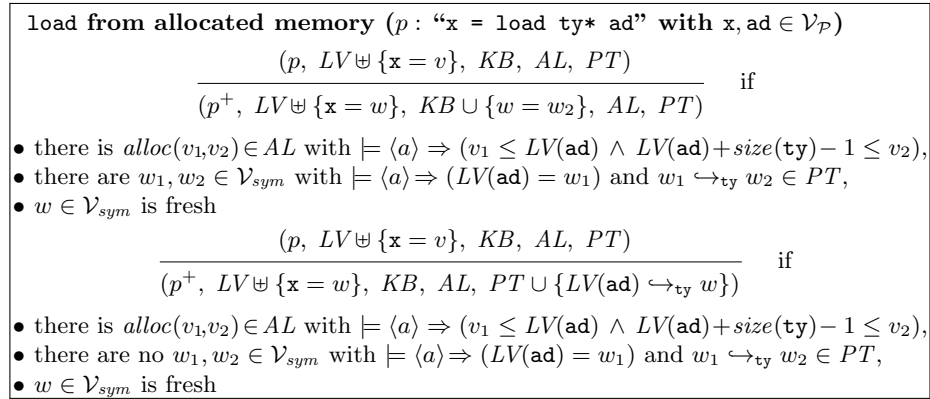


Fig. 1. Symbolic execution graph for `strlen`

$LV(\text{ad}) + \text{size}(\text{ty}) - 1 \leq v_2$ ) is valid. Then, we reach a new abstract state where the previous position  $p = (\mathbf{b}_{prev}, \mathbf{b}, i)$  is updated to the position  $p^+ = (\mathbf{b}_{prev}, \mathbf{b}, i + 1)$  of the next instruction in the same basic block, and we set  $LV(\mathbf{x}) = w$  for a fresh  $w \in \mathcal{V}_{sym}$ . If we already know the value at the address  $\text{ad}$  (i.e., if there are  $w_1, w_2 \in \mathcal{V}_{sym}$  with  $\models \langle a \rangle \Rightarrow (LV(\text{ad}) = w_1)$  and  $w_1 \hookrightarrow_{\text{ty}} w_2 \in PT$ ) then we add  $w = w_2$  to  $KB$ . Otherwise, we add  $LV(\text{ad}) \hookrightarrow_{\text{ty}} w$  to  $PT$ . We used this rule to obtain  $B$  from  $A$  in Fig. 1. In a similar way, one can also formulate a rule for `store` instructions that store a value at some address in the memory (cf. [26]).



If `load` accesses an address that was not allocated, then memory safety is violated and we reach the *ERR* state.

<p style="text-align: center;"><b>load from unallocated memory</b> (<math>p</math>: “<code>x = load ty* ad</code>” with <math>x, ad \in \mathcal{V}_P</math>)</p> $\frac{(p, LV, KB, AL, PT)}{ERR} \quad \text{if}$ <p>there is no <math>alloc(v_1, v_2) \in AL</math> with <math>\models \langle a \rangle \Rightarrow (v_1 \leq LV(ad) \wedge LV(ad) + size(ty) - 1 \leq v_2)</math></p>
--

The instructions `icmp` and `br` in `strlen`’s entry block check if the first character `c0` is 0. In that case, we have reached the end of the string and jump to the block `done`. So for “`x = icmp eq ty t1, t2`”, we check if the state contains enough information to decide whether the values  $t_1$  and  $t_2$  of type `ty` are equal. In that case, the value 1 resp. 0 (i.e., *true* resp. *false*) is assigned to `x`.<sup>11</sup>

<p style="text-align: center;"><b>icmp</b> (<math>p</math>: “<code>x = icmp eq ty t<sub>1</sub>, t<sub>2</sub></code>” with <math>x \in \mathcal{V}_P</math> and <math>t_1, t_2 \in \mathcal{V}_P \cup \mathbb{Z}</math>)</p> $\frac{(p, LV \uplus \{x = v\}, KB, AL, PT)}{(p^+, LV \uplus \{x = w\}, KB \cup \{w = 1\}, AL, PT)} \quad \begin{array}{l} \text{if } \models \langle a \rangle \Rightarrow (LV(t_1) = LV(t_2)) \\ \text{and } w \in \mathcal{V}_{sym} \text{ is fresh} \end{array}$ $\frac{(p, LV \uplus \{x = v\}, KB, AL, PT)}{(p^+, LV \uplus \{x = w\}, KB \cup \{w = 0\}, AL, PT)} \quad \begin{array}{l} \text{if } \models \langle a \rangle \Rightarrow (LV(t_1) \neq LV(t_2)) \\ \text{and } w \in \mathcal{V}_{sym} \text{ is fresh} \end{array}$
--

The previous rule is only applicable if *KB* contains enough information to evaluate the condition. Otherwise, a case analysis needs to be performed, i.e., one has to *refine* the abstract state by extending its knowledge base. This is done by the following rule which transforms an abstract state into *two* new ones.<sup>12</sup>

<p style="text-align: center;"><b>refining abstract states</b> (<math>p</math>: “<code>x = icmp eq ty t<sub>1</sub>, t<sub>2</sub></code>”, <math>x \in \mathcal{V}_P</math>, <math>t_1, t_2 \in \mathcal{V}_P \cup \mathbb{Z}</math>)</p> $\frac{(p, LV, KB, AL, PT)}{(p, LV, KB \cup \{\varphi\}, AL, PT) \mid (p, LV, KB \cup \{\neg\varphi\}, AL, PT)}$ <p>if <math>\varphi</math> is <math>LV(t_1) = LV(t_2)</math> and both <math>\not\models \langle a \rangle \Rightarrow \varphi</math> and <math>\not\models \langle a \rangle \Rightarrow \neg\varphi</math></p>
---

For example, in state *B* of Fig. 1, we evaluate “`c0zero = icmp eq i8 c0, 0`”, i.e., we check whether the first character `c0` of the string `str` is 0. Since this cannot be inferred from *B*’s knowledge base, we refine *B* to the successor states *C* and *D* and call the edges from *B* to *C* and *D* *refinement edges*. In *D*, we have `c0 = v1` and  $v_1 \neq 0$ . Thus, the `icmp`-rule yields *E* where `c0zero = v2` and  $v_2 = 0$ . We do not display the successors of *C* that lead to a program end.

The conditional branching instruction `br` is very similar to `icmp`. To evaluate “`br i1 t, label b1, label b2`”, one has to check whether the current state contains enough information to conclude that  $t$  is 1 (i.e., *true*) or 0 (i.e., *false*). Then the evaluation continues with block `b1` resp. `b2`. This rule allows us to create the successor *F* of *E*, where we jump to the block `loop`.

<sup>11</sup> Other integer comparisons (for  $<$ ,  $\leq$ , ...) are handled analogously.

<sup>12</sup> Analogous refinement rules can also be used for other conditional LLVM instructions.



$$\begin{array}{c}
\text{br } (p : \text{“br i1 } t, \text{label } \mathbf{b}_1, \text{label } \mathbf{b}_2\text{” with } t \in \mathcal{V}_{\mathcal{P}} \cup \{0, 1\} \text{ and } \mathbf{b}_1, \mathbf{b}_2 \in \text{Blks}) \\
\frac{(p, LV, KB, AL, PT)}{((\mathbf{b}, \mathbf{b}_1, 0), LV, KB, AL, PT)} \quad \text{if } p = (\mathbf{b}_{prev}, \mathbf{b}, i) \text{ and } \models \langle a \rangle \Rightarrow (LV(t) = 1) \\
\frac{(p, LV, KB, AL, PT)}{((\mathbf{b}, \mathbf{b}_2, 0), LV, KB, AL, PT)} \quad \text{if } p = (\mathbf{b}_{prev}, \mathbf{b}, i) \text{ and } \models \langle a \rangle \Rightarrow (LV(t) = 0)
\end{array}$$

Next, we have to evaluate a **phi** instruction. These instructions are needed due to the static single assignment form of LLVM. Here, “ $\mathbf{x} = \text{phi ty } [t_1, \mathbf{b}_1], \dots, [t_n, \mathbf{b}_n]$ ” means that if the previous block was  $\mathbf{b}_j$ , then the value  $t_j$  is assigned to  $\mathbf{x}$ . All  $t_1, \dots, t_n$  must have type  $\text{ty}$ . Since we reached state  $F$  in Fig. 1 after evaluating the **entry** block, we obtain the state  $G$  with  $\text{olds} = v_3$  and  $v_3 = u_{\text{str}}$ .

$$\begin{array}{c}
\text{phi } (p : \text{“x = phi ty } [t_1, \mathbf{b}_1], \dots, [t_n, \mathbf{b}_n]\text{” with } \mathbf{x} \in \mathcal{V}_{\mathcal{P}}, t_i \in \mathcal{V}_{\mathcal{P}} \cup \mathbb{Z}, \mathbf{b}_i \in \text{Blks}) \\
\frac{(p, LV \uplus \{\mathbf{x} = v\}, KB, AL, PT)}{(p^+, LV \uplus \{\mathbf{x} = w\}, KB \cup \{w = LV(t_j)\}, AL, PT)} \quad \text{if } p = (\mathbf{b}_j, \mathbf{b}, k) \text{ and } \\
\quad w \in \mathcal{V}_{\text{sym}} \text{ is fresh}
\end{array}$$

The **strlen** function traverses the string using a pointer  $\mathbf{s}$  that is increased in each iteration. The **loop** terminates, since eventually  $\mathbf{s}$  reaches the last memory cell of the string (containing 0). Then one jumps to **done**, converts the pointers  $\mathbf{s}$  and **str** to integers, and returns their difference. To perform the required pointer arithmetic, “ $\text{ad}_2 = \text{getelementptr ty* ad}_1, \text{in } t$ ” increases  $\text{ad}_1$  by the size of  $t$  elements of type  $\text{ty}$  (i.e., by  $\text{size}(\text{ty}) \cdot t$ ) and assigns this address to  $\text{ad}_2$ .<sup>13</sup>

$$\begin{array}{c}
\text{getelementptr } (p : \text{“ad}_2 = \text{getelementptr ty* ad}_1, \text{in } t\text{”, } \text{ad}_1, \text{ad}_2 \in \mathcal{V}_{\mathcal{P}}, t \in \mathcal{V}_{\mathcal{P}} \cup \mathbb{Z}) \\
\frac{(p, LV \uplus \{\text{ad}_2 = v\}, KB, AL, PT)}{(p^+, LV \uplus \{\text{ad}_2 = w\}, KB \cup \{w = LV(\text{ad}_1) + \text{size}(\text{ty}) \cdot LV(t)\}, AL, PT)} \quad w \in \mathcal{V}_{\text{sym}} \text{ fresh}
\end{array}$$

In Fig. 1, this rule is used for the step from  $G$  to  $H$ , where  $LV$  and  $KB$  now imply  $\mathbf{s} = \text{str} + 1$ . In the step to  $I$ , the character at address  $\mathbf{s}$  is loaded to  $\mathbf{c}$ . To ensure memory safety, the **load**-rule checks that  $\mathbf{s}$  is in an allocated part of the memory (i.e., that  $u_{\text{str}} \leq u_{\text{str}} + 1 \leq v_{\text{end}}$ ). This holds because  $\langle H \rangle$  implies  $u_{\text{str}} \leq v_{\text{end}}$  and  $u_{\text{str}} \neq v_{\text{end}}$  (as  $u_{\text{str}} \hookrightarrow v_1, v_{\text{end}} \hookrightarrow 0 \in PT$  and  $v_1 \neq 0 \in KB$ ). Finally, we check whether  $\mathbf{c}$  is 0. We again perform a refinement which yields the states  $J$  and  $K$ . State  $K$  corresponds to the case  $\mathbf{c} \neq 0$  and thus, we obtain  $\text{czero} = 0$  in  $L$  and branch back to instruction 0 of the **loop** block in state  $M$ .

### 2.3 Generalizing Abstract States

After reaching  $M$ , one unfolds the loop once more until one reaches a state  $\tilde{M}$  at position  $(\text{loop}, \text{loop}, 0)$  again, analogous to the first iteration. To obtain *finite* symbolic execution graphs, we *generalize* our states whenever an evaluation visits a program position twice. Thus, we have to find a state that is more general than  $M = (p, LV_M, KB_M, AL, PT_M)$  and  $\tilde{M} = (p, LV_{\tilde{M}}, KB_{\tilde{M}}, AL, PT_{\tilde{M}})$ . For readability, we again write “ $\hookrightarrow$ ” instead of “ $\hookrightarrow_{i8}$ ”. Then  $p = (\text{loop}, \text{loop}, 0)$  and

<sup>13</sup> Since we do not consider the handling of data structures in this paper, we do not regard **getelementptr** instructions with more than two parameters.

$$\begin{aligned}
AL &= \{alloc(u_{\mathbf{str}}, v_{end})\} \\
LV_M &= \{\mathbf{str} = u_{\mathbf{str}}, \mathbf{c} = v_5, \mathbf{s} = v_4, \mathbf{olds} = v_3, \dots\} \\
LV_{\widetilde{M}} &= \{\mathbf{str} = u_{\mathbf{str}}, \mathbf{c} = \widetilde{v}_5, \mathbf{s} = \widetilde{v}_4, \mathbf{olds} = \widetilde{v}_3, \dots\} \\
PT_M &= \{u_{\mathbf{str}} \leftrightarrow v_1, v_4 \leftrightarrow v_5, v_{end} \leftrightarrow z\} \\
PT_{\widetilde{M}} &= \{u_{\mathbf{str}} \leftrightarrow v_1, v_4 \leftrightarrow v_5, \widetilde{v}_4 \leftrightarrow \widetilde{v}_5, v_{end} \leftrightarrow z\} \\
KB_M &= \{v_5 \neq 0, v_4 = v_3 + 1, v_3 = u_{\mathbf{str}}, v_1 \neq 0, z = 0, \dots\} \\
KB_{\widetilde{M}} &= \{\widetilde{v}_5 \neq 0, \widetilde{v}_4 = \widetilde{v}_3 + 1, \widetilde{v}_3 = v_4, v_4 = v_3 + 1, v_3 = u_{\mathbf{str}}, v_1 \neq 0, z = 0, \dots\}.
\end{aligned}$$

Our aim is to construct a new state  $N$  that is more general than  $M$  and  $\widetilde{M}$ , but contains enough information for the remaining proof. We now present our heuristic for *merging* states that is used as the basis for our implementation.

To merge  $M$  and  $\widetilde{M}$ , we keep those constraints of  $M$  that also hold in  $\widetilde{M}$ . To this end, we proceed in two steps. First, we create a new state  $N = (p, LV_N, KB_N, AL_N, PT_N)$  using fresh symbolic variables  $v_x$  for all  $\mathbf{x} \in \mathcal{V}_{\mathcal{P}}$  and define

$$LV_N = \{\mathbf{str} = v_{\mathbf{str}}, \mathbf{c} = v_c, \mathbf{s} = v_s, \mathbf{olds} = v_{\mathbf{olds}}, \dots\}.$$

Matching  $N$ 's fresh variables to the variables in  $M$  and  $\widetilde{M}$  yields mappings with  $\mu_M(v_{\mathbf{str}}) = u_{\mathbf{str}}, \mu_M(v_c) = v_5, \mu_M(v_s) = v_4, \mu_M(v_{\mathbf{olds}}) = v_3$ , and  $\mu_{\widetilde{M}}(v_{\mathbf{str}}) = u_{\mathbf{str}}, \mu_{\widetilde{M}}(v_c) = \widetilde{v}_5, \mu_{\widetilde{M}}(v_s) = \widetilde{v}_4, \mu_{\widetilde{M}}(v_{\mathbf{olds}}) = \widetilde{v}_3$ . By injectivity of  $LV_M$ , we can also define a pseudo-inverse of  $\mu_M$  that maps  $M$ 's variables to  $N$  by setting  $\mu_M^{-1}(LV_M(\mathbf{x})) = v_x$  for  $\mathbf{x} \in \mathcal{V}_{\mathcal{P}}$  and  $\mu_M^{-1}(v) = v$  for all other  $v \in \mathcal{V}_{sym}$  ( $\mu_{\widetilde{M}}^{-1}$  works analogously).

In a second step, we use these mappings to check which constraints of  $M$  also hold in  $\widetilde{M}$ . So we set  $AL_N = \mu_M^{-1}(AL) \cap \mu_{\widetilde{M}}^{-1}(AL) = \{alloc(v_{\mathbf{str}}, v_{end})\}$  and

$$\begin{aligned}
PT_N &= \mu_M^{-1}(PT_M) \cap \mu_{\widetilde{M}}^{-1}(PT_{\widetilde{M}}) \\
&= \{v_{\mathbf{str}} \leftrightarrow v_1, v_s \leftrightarrow v_c, v_{end} \leftrightarrow z\} \cap \{v_{\mathbf{str}} \leftrightarrow v_1, v_4 \leftrightarrow v_5, v_s \leftrightarrow v_c, v_{end} \leftrightarrow z\} \\
&= \{v_{\mathbf{str}} \leftrightarrow v_1, v_s \leftrightarrow v_c, v_{end} \leftrightarrow z\}.
\end{aligned}$$

It remains to construct  $KB_N$ . We have  $v_3 = u_{\mathbf{str}}$  (“olds = str”) in  $\langle M \rangle$ , but  $\widetilde{v}_3 = v_4, v_4 = v_3 + 1, v_3 = u_{\mathbf{str}}$  (“olds = str + 1”) in  $\langle \widetilde{M} \rangle$ . To keep as much information as possible in such cases, we rewrite equations to inequations before performing the generalization. For this, let  $\langle\langle M \rangle\rangle$  result from extending  $\langle M \rangle$  by  $t_1 \geq t_2$  and  $t_1 \leq t_2$  for any equation  $t_1 = t_2 \in \langle M \rangle$ . So in our example, we obtain  $v_3 \geq u_{\mathbf{str}} \in \langle\langle M \rangle\rangle$  (“olds  $\geq$  str”). Moreover, for any  $t_1 \neq t_2 \in \langle M \rangle$ , we check whether  $\langle M \rangle$  implies  $t_1 > t_2$  or  $t_1 < t_2$ , and add the respective inequation to  $\langle\langle M \rangle\rangle$ . In this way, one can express sequences of inequations  $t_1 \neq t_2, t_1 + 1 \neq t_2, \dots, t_1 + n \neq t_2$  (where  $t_1 \leq t_2$ ) by a single inequation  $t_1 + n < t_2$ , which is needed for suitable generalizations afterwards. We use this to derive  $v_4 < v_{end} \in \langle\langle M \rangle\rangle$  (“s < v\_end”) from  $v_4 = v_3 + 1, v_3 = u_{\mathbf{str}}, u_{\mathbf{str}} \leq v_{end}, u_{\mathbf{str}} \neq v_{end}, v_4 \neq v_{end} \in \langle M \rangle$ .

We then let  $KB_N$  consist of all formulas  $\varphi$  from  $\langle\langle M \rangle\rangle$  that are also implied by  $\langle \widetilde{M} \rangle$ , again translating variable names using  $\mu_M^{-1}$  and  $\mu_{\widetilde{M}}^{-1}$ . Thus, we have

$$\begin{aligned}
\langle\langle M \rangle\rangle &= \{v_5 \neq 0, v_4 = v_3 + 1, v_3 = u_{\mathbf{str}}, v_3 \geq u_{\mathbf{str}}, v_4 < v_{end}, \dots\} \\
\mu_M^{-1}(\langle\langle M \rangle\rangle) &= \{v_c \neq 0, v_s = v_{\mathbf{olds}} + 1, v_{\mathbf{olds}} = v_{\mathbf{str}}, v_{\mathbf{olds}} \geq v_{\mathbf{str}}, v_s < v_{end}, \dots\} \\
\mu_{\widetilde{M}}^{-1}(\langle\langle \widetilde{M} \rangle\rangle) &= \{v_c \neq 0, v_s = v_{\mathbf{olds}} + 1, v_{\mathbf{olds}} = v_4, v_4 = v_3 + 1, v_3 = u_{\mathbf{str}}, v_s < v_{end}, \dots\} \\
KB_N &= \{v_c \neq 0, v_s = v_{\mathbf{olds}} + 1, v_{\mathbf{olds}} \geq v_{\mathbf{str}}, v_s < v_{end}, \dots\}.
\end{aligned}$$

**Definition 3 (Merging States).** Let  $a = (p, LV_a, KB_a, AL_a, PT_a)$  and  $b = (p, LV_b, KB_b, AL_b, PT_b)$  be abstract states. Then  $c = (p, LV_c, KB_c, AL_c, PT_c)$  results from merging the states  $a$  and  $b$  if

- $LV_c = \{\mathbf{x} = v_{\mathbf{x}} \mid \mathbf{x} \in \mathcal{V}_{\mathcal{P}}\}$  for fresh pairwise different symbolic variables  $v_{\mathbf{x}}$ . Moreover, we define  $\mu_a(v_{\mathbf{x}}) = LV_a(\mathbf{x})$  and  $\mu_b(v_{\mathbf{x}}) = LV_b(\mathbf{x})$  for all  $\mathbf{x} \in \mathcal{V}_{\mathcal{P}}$  and let  $\mu_a$  and  $\mu_b$  be the identity on all remaining variables from  $\mathcal{V}_{sym}$ .
- $AL_c = \mu_a^{-1}(AL_a) \cap \mu_b^{-1}(AL_b)$  and  $PT_c = \mu_a^{-1}(PT_a) \cap \mu_b^{-1}(PT_b)$ . Here, the “inverse” of the instantiation  $\mu_a$  is defined as  $\mu_a^{-1}(v) = v_{\mathbf{x}}$  if  $v = LV_a(\mathbf{x})$  and  $\mu_a^{-1}(v) = v$  for all other  $v \in \mathcal{V}_{sym}$  ( $\mu_b^{-1}$  is defined analogously).
- $KB_c = \{ \varphi \in \mu_a^{-1}(\langle\langle a \rangle\rangle) \mid \models \mu_b^{-1}(\langle\langle b \rangle\rangle) \Rightarrow \varphi \}$ , where
 
$$\begin{aligned} \langle\langle a \rangle\rangle &= \langle a \rangle \cup \{ t_1 \geq t_2, t_1 \leq t_2 \mid t_1 = t_2 \in \langle a \rangle \} \\ &\cup \{ t_1 > t_2 \mid t_1 \neq t_2 \in \langle a \rangle, \models \langle a \rangle \Rightarrow t_1 > t_2 \} \\ &\cup \{ t_1 < t_2 \mid t_1 \neq t_2 \in \langle a \rangle, \models \langle a \rangle \Rightarrow t_1 < t_2 \}. \end{aligned}$$

In Fig. 1, we do not show the second loop unfolding from  $M$  to  $\widetilde{M}$ , and directly draw a *generalization edge* from  $M$  to  $N$ , depicted by a dashed arrow. Such an edge expresses that all concrete states represented by  $M$  are also represented by  $N$ . Semantically, a state  $\bar{a}$  is a generalization of a state  $a$  iff  $\models \langle a \rangle_{SL} \Rightarrow \mu(\langle \bar{a} \rangle_{SL})$  for some instantiation  $\mu$ . To automate our procedure, we define a weaker relationship between  $a$  and  $\bar{a}$ . We say that  $\bar{a} = (p, \overline{LV}, \overline{KB}, \overline{AL}, \overline{PT})$  is a *generalization* of  $a = (p, LV, KB, AL, PT)$  with the instantiation  $\mu$  whenever the conditions (b)-(e) of the following rule are satisfied.

<p><b>generalization with <math>\mu</math></b></p> $\frac{(p, LV, KB, AL, PT)}{(p, \overline{LV}, \overline{KB}, \overline{AL}, \overline{PT})} \text{ if}$ <p>(a) <math>a</math> has an incoming evaluation edge,<sup>14</sup>          (b) <math>LV(\mathbf{x}) = \mu(\overline{LV}(\mathbf{x}))</math> for all <math>\mathbf{x} \in \mathcal{V}_{\mathcal{P}}</math>,          (c) <math>\models \langle a \rangle \Rightarrow \mu(\overline{KB})</math>,          (d) if <math>alloc(v_1, v_2) \in \overline{AL}</math>, then <math>alloc(\mu(v_1), \mu(v_2)) \in AL</math>,          (e) if <math>(v_1 \hookrightarrow_{ty} v_2) \in \overline{PT}</math>, then <math>(\mu(v_1) \hookrightarrow_{ty} \mu(v_2)) \in PT</math></p>
---

Clearly, then we indeed have  $\models \langle a \rangle_{SL} \Rightarrow \mu(\langle \bar{a} \rangle_{SL})$ . Condition (a) is needed to avoid cycles of refinement and generalization steps in the symbolic execution graph, which would not correspond to any computation.

Of course, many approaches are possible to compute such generalizations (or “widenings”). Thm. 4 shows that the merging heuristic from Def. 3 satisfies the conditions of the generalization rule. Thus, since  $N$  results from merging  $M$  and  $\widetilde{M}$ , it is indeed a *generalization* of  $M$ . Thm. 4 also shows that if one uses the merging heuristic to compute generalizations, then the construction of symbolic execution graphs always terminates when applying the following strategy:

- If there is a path from a state  $a$  to a state  $b$ , where  $a$  and  $b$  are at the same program position, where  $b$  has an incoming evaluation edge, and where  $a$  has no incoming refinement edge, then we check whether  $a$  is a generalization of

<sup>14</sup> *Evaluation edges* are edges that are not refinement or generalization edges.

$b$  (i.e., whether the corresponding conditions of the generalization rule are satisfied). In that case, we draw a generalization edge from  $b$  to  $a$ .

- Otherwise, remove  $a$ 's children, and add a generalization edge from  $a$  to the merging  $c$  of  $a$  and  $b$ . If  $a$  already had an incoming generalization edge from some state  $q$ , then remove  $a$  and add a generalization edge from  $q$  to  $c$  instead.

**Theorem 4 (Soundness and Termination of Merging).** *Let  $c$  result from merging the states  $a$  and  $b$  as in Def. 3. Then  $c$  is a generalization of  $a$  and  $b$  with the instantiations  $\mu_a$  and  $\mu_b$ , respectively. Moreover, if  $a$  is not already a generalization of  $b$ , then  $|\langle\langle c \rangle\rangle| + |AL_c| + |PT_c| < |\langle\langle a \rangle\rangle| + |AL_a| + |PT_a|$ . Here, for any conjunction  $\varphi$ , let  $|\varphi|$  denote the number of its conjuncts. Thus, the above strategy to construct symbolic execution graphs always terminates.<sup>15</sup>*

In our example, we continue symbolic execution in state  $N$ . Similar to the execution from  $F$  to  $M$ , after 6 steps another state  $P$  at position  $(\text{loop}, \text{loop}, 0)$  is reached. In Fig. 1, dotted arrows abbreviate several evaluation steps. As  $N$  is again a generalization of  $P$  using an instantiation  $\mu$  with  $\mu(v_c) = w_c$ ,  $\mu(v_s) = w_s$ , and  $\mu(v_{\text{oids}}) = w_{\text{oids}}$ , we draw a generalization edge from  $P$  to  $N$ . The construction of the symbolic execution graph is finished as soon as all its leaves correspond to `ret` instructions (for “return”).

Based on this construction, we now connect the symbolic execution graph to memory safety of the input program. We say that a concrete LLVM state  $(p, s, m)$  is *represented* by the symbolic execution graph iff the graph contains an abstract state  $a$  at position  $p$  where  $(s, m) \models \sigma(\langle a \rangle_{SL})$  for some concrete instantiation  $\sigma$ .

**Theorem 5 (Memory Safety of LLVM Programs).** *Let  $\mathcal{P}$  be an LLVM program with a symbolic execution graph  $\mathcal{G}$ . If  $\mathcal{G}$  does not contain the abstract state  $ERR$ , then  $\mathcal{P}$  is memory safe for all LLVM states represented by  $\mathcal{G}$ .*

### 3 From Symbolic Execution Graphs to Integer Systems

To prove termination of the input program, we extract an *integer transition system* (ITS) from the symbolic execution graph and then use existing tools to prove its termination. The extraction step essentially restricts the information in abstract states to the integer constraints on symbolic variables. This conversion of memory-based arguments into integer arguments often suffices for the termination proof. The reason for considering only  $\mathcal{V}_{sym}$  instead of  $\mathcal{V}_{\mathcal{P}}$  is that the conditions in the abstract states only concern the symbolic variables and therefore, these are usually the essential variables for proving termination.

For example, termination of `strlen` is proved by showing that the pointer `s` is increased as long as it is smaller than  $v_{end}$ , the symbolic end of the input string. In Fig. 1, this is explicit since  $v_s < v_{end}$  is an invariant that holds in all states represented by  $N$ . Each iteration of the cycle increases the value of  $v_s$ .

Formally, *ITSs* are graphs whose nodes are abstract states and whose edges are *transitions*. For any abstract state  $a$ , let  $\mathcal{V}(a)$  denote the symbolic variables occurring in  $a$ . Let  $\mathcal{V} \subseteq \mathcal{V}_{sym}$  be the finite set of all symbolic variables occurring

<sup>15</sup> The proofs for all theorems can be found in [26].

in states of the symbolic execution graph. A *transition* is a tuple  $(a, CON, \bar{a})$  where  $a, \bar{a}$  are abstract states and the *condition*  $CON \subseteq QF\_IA(\mathcal{V} \uplus \mathcal{V}')$  is a set of pure quantifier-free formulas over the variables  $\mathcal{V} \uplus \mathcal{V}'$ . Here,  $\mathcal{V}' = \{v' \mid v \in \mathcal{V}\}$  represents the values of the variables *after* the transition. An *ITS state*  $(a, \sigma)$  consists of an abstract state  $a$  and a concrete instantiation  $\sigma : \mathcal{V} \rightarrow \mathbb{Z}$ . For any such  $\sigma$ , let  $\sigma' : \mathcal{V}' \rightarrow \mathbb{Z}$  with  $\sigma'(v') = \sigma(v)$ . Given an ITS  $\mathcal{I}$ ,  $(a, \sigma)$  *evaluates* to  $(\bar{a}, \bar{\sigma})$  (denoted “ $(a, \sigma) \rightarrow_{\mathcal{I}} (\bar{a}, \bar{\sigma})$ ”) iff  $\mathcal{I}$  has a transition  $(a, CON, \bar{a})$  with  $\models (\sigma \cup \sigma')(CON)$ . Here, we have  $(\sigma \cup \sigma')(v) = \sigma(v)$  and  $(\sigma \cup \sigma')(v') = \sigma'(v') = \bar{\sigma}(v')$  for all  $v \in \mathcal{V}$ . An ITS  $\mathcal{I}$  is *terminating* iff  $\rightarrow_{\mathcal{I}}$  is well founded.<sup>16</sup>

We convert symbolic execution graphs to ITSs by transforming every edge into a transition. If there is a generalization edge from  $a$  to  $\bar{a}$  with an instantiation  $\mu$ , then the new value of any  $v \in \mathcal{V}(\bar{a})$  in  $\bar{a}$  is  $\mu(v)$ . Hence, we create the transition  $(a, \langle a \rangle \cup \{v' = \mu(v) \mid v \in \mathcal{V}(\bar{a})\}, \bar{a})$ .<sup>17</sup> So for the edge from  $P$  to  $N$  in Fig. 1, we obtain the condition  $\{w_s = w_{\text{old}_s} + 1, w_{\text{old}_s} = v_s, v_s < v_{\text{end}}, v'_{\text{str}} = v_{\text{str}}, v'_{\text{end}} = v_{\text{end}}, v'_c = w_c, v'_s = w_s, \dots\}$ . This can be simplified to  $\{v_s < v_{\text{end}}, v'_{\text{end}} = v_{\text{end}}, v'_s = v_s + 1, \dots\}$ .

An evaluation or refinement edge from  $a$  to  $\bar{a}$  does not change the variables of  $\mathcal{V}(a)$ . Thus, we construct the transition  $(a, \langle a \rangle \cup \{v' = v \mid v \in \mathcal{V}(a)\}, \bar{a})$ .

So in the ITS resulting from Fig. 1, the condition of the transition from  $A$  to  $B$  contains  $\{v'_{\text{end}} = v_{\text{end}}\} \cup \{u'_x = u_x \mid x \in \mathcal{V}_{\mathcal{P}}\}$ . The condition for the transition from  $B$  to  $D$  is the same, but extended by  $v'_1 = v_1$ . Hence, in the transition from  $A$  to  $B$ , the value of  $v_1$  can change arbitrarily (since  $v_1 \notin \mathcal{V}(A)$ ), but in the transition from  $B$  to  $D$ , the value of  $v_1$  must remain the same.

**Definition 6 (ITS from Symbolic Execution Graph).** *Let  $\mathcal{G}$  be a symbolic execution graph. Then the corresponding integer transition system  $\mathcal{I}_{\mathcal{G}}$  has one transition for each edge in  $\mathcal{G}$ :*

- *If the edge from  $a$  to  $\bar{a}$  is not a generalization edge, then  $\mathcal{I}_{\mathcal{G}}$  has a transition from  $a$  to  $\bar{a}$  with the condition  $\langle a \rangle \cup \{v' = v \mid v \in \mathcal{V}(a)\}$ .*
- *If there is a generalization edge from  $a$  to  $\bar{a}$  with the instantiation  $\mu$ , then  $\mathcal{I}_{\mathcal{G}}$  has a transition from  $a$  to  $\bar{a}$  with the condition  $\langle a \rangle \cup \{v' = \mu(v) \mid v \in \mathcal{V}(\bar{a})\}$ .*

From the non-generalization edges on the path from  $N$  to  $P$  in Fig. 1, we obtain transitions whose conditions contain  $v'_{\text{end}} = v_{\text{end}}$  and  $v'_s = v_s$ . So  $v_s$  is increased by 1 in the transition from  $P$  to  $N$  and it remains the same in all other transitions of the graph’s only cycle. Since the transition from  $P$  to  $N$  is only executed as long as  $v_s < v_{\text{end}}$  holds (where  $v_{\text{end}}$  is not changed by any transition), termination of the resulting ITS can easily be proved automatically.

The following theorem shows the soundness of our approach.

**Theorem 7 (Termination of LLVM Programs).** *Let  $\mathcal{P}$  be an LLVM program with a symbolic execution graph  $\mathcal{G}$  that does not contain the state  $ERR$ . If  $\mathcal{I}_{\mathcal{G}}$  is terminating, then  $\mathcal{P}$  is also terminating for all LLVM states represented by  $\mathcal{G}$ .*

<sup>16</sup> For programs starting in states represented by an abstract state  $a_0$ , it would suffice to prove termination of all  $\rightarrow_{\mathcal{I}}$ -evaluations starting in ITS states of the form  $(a_0, \sigma)$ .

<sup>17</sup> In the transition, we do not impose the additional constraints of  $\langle \bar{a} \rangle$  on the post-variables  $\mathcal{V}'$ , since they are checked anyway in the next transition which starts in  $\bar{a}$ .

## 4 Related Work, Experiments, and Conclusion

We developed a new approach to prove memory safety and termination of C (resp. LLVM) programs with explicit pointer arithmetic and memory access. It relies on a representation of abstract program states which allows an easy automation of the rules for symbolic execution (by standard SMT solving). Moreover, this representation is suitable for generalizing abstract states and for generating integer transition systems. In this way, LLVM programs are translated fully automatically into ITSs amenable to automated termination analysis.

Previous methods and tools for termination analysis of imperative programs (e.g., AProVE [4, 5], ARMC [24], COSTA [1], Cyclist [7], FuncTion [29], Julia [25], KITTeL [12], LoopFrog [28], TAN [16], TRex [14], T2 [6], Ultimate [15], ...) either do not handle the heap at all, or support dynamic data structures by an abstraction to integers (e.g., to represent sizes or lengths) or to terms (representing finite unravelings). However, most tools fail when the control flow depends on explicit pointer arithmetic and on detailed information about the contents of addresses. While the general methodology of our approach was inspired by our previous work on termination of Java [4, 5], in the current paper we lift such techniques to prove termination and memory safety of programs with explicit pointer arithmetic. This requires a fundamentally new approach, since pointer arithmetic and memory allocation cannot be expressed in the Java-based techniques of [4, 5].

We implemented our technique in the termination prover AProVE using the SMT solvers Yices [11] and Z3 [20] in the back-end. A preliminary version of our implementation participated very successfully in the *International Competition on Software Verification (SV-COMP)* [27] at TACAS, which featured a category for termination of C programs for the first time in 2014. To evaluate AProVE’s power, we performed experiments on a collection of 208 C programs from several sources, including the *SV-COMP 2014* termination category and standard string algorithms from [30] and the OpenBSD C library [23]. Of these 208 programs, 129 use pointers and 79 only operate on integers.

To prove termination of low-level C programs, one also has to ensure their memory safety. While there exist several tools to prove memory safety of C programs, many of them do not handle explicit byte-accurate pointer arithmetic (e.g., Thor [19] or SLayer [3]) or require the user to provide the needed loop invariants (as in the Jessie plug-in of Frama-C [21]). In contrast, our approach can prove memory safety of such algorithms fully automatically. Although our approach is targeted toward termination and only analyzes memory safety as a prerequisite for termination, it turned out that on our collection, AProVE is more powerful than the leading publicly available tools for proving memory safety. To this end, we compared AProVE with the tools CPAchecker [18] and Predator [10] which reached the first and the third place in the category for *memory safety* at *SV-COMP 2014*.<sup>18</sup> For the 129 pointer programs in our collection, AProVE can show memory safety for 102 examples, whereas CPAchecker resp. Predator prove

<sup>18</sup> The second place in this category was reached by the bounded model checker LLBMC [13]. However, in general such tools only disprove, but cannot verify memory safety.

memory safety for 77 resp. 79 examples (see [2] for details).

To evaluate the power of our approach for proving termination, we compared AProVE to the other tools from the termination category of *SV-COMP 2014*. In addition, we included the termination analyzer KITTeL [12] in our evaluation,

which operates on LLVM as well. On the side, we show the performance of the tools on integer and pointer programs when using a time limit of 300 seconds for

	79 integer programs					129 pointer programs				
	T	N	F	TO	RT	T	N	F	TO	RT
AProVE	67	0	11	1	19.6	91	0	19	19	58.6
FuncTion	11	0	66	2	23.1	-	-	-	-	-
KITTeL	58	0	12	9	0.2	9	0	1	119	0.2
T2	55	0	23	1	1.8	6	0	123	0	3.6
TAN	31	0	37	11	2.4	3	0	124	2	10.6
Ultimate	57	4	12	6	3.2	-	-	-	-	-

each example. Here, we used an Intel Core i7-950 processor and 6 GB of memory. “**T**” gives the number of examples where termination could be proved, “**N**” is the number of examples where non-termination could be shown, “**F**” states how often the tool failed in less than 300 seconds, “**TO**” gives the number of time-outs (i.e., examples for which the tool took longer than 300 seconds), and “**RT**” is the average run time in seconds for those examples where the tool proved termination or non-termination. For pointer programs, we omitted the results for those tools that were not able to prove termination of any examples.

Most other termination provers ignore the problem of memory safety and just prove termination under the *assumption* that the program is memory safe. So they may also return “Yes” for memory unsafe programs and may treat read accesses to the heap as non-deterministic input. Since AProVE constructs symbolic execution graphs to prove memory safety and to infer suitable invariants needed for termination proofs, its runtime is often higher than that of other tools. On the other hand, the table shows that our approach is slightly more powerful than the other tools for integer programs (i.e., our graph-based technique is also suitable for programs on integers) and it is clearly the most powerful one for pointer programs. The reason is due to our novel representation of the memory which handles pointer arithmetic and keeps information about the contents of addresses. For details on our experiments and to access our implementation in AProVE via a web interface, we refer to [2]. In future work, we plan to extend our approach to recursive programs and to inductive data structures defined via `struct` (e.g., by integrating existing shape analyses based on separation logic).

**Acknowledgments.** We are grateful to the developers of the other tools for termination or memory safety [6, 10, 12, 15, 16, 18, 29] for their help with the experiments.

## References

1. Albert, E., Arenas, P., Codish, M., Genaim, S., Puebla, G., Zanardini, D.: Termination analysis of Java Bytecode. In: Proc. FMOODS ’08
2. AProVE: <http://aprove.informatik.rwth-aachen.de/eval/Pointer/>
3. Berdine, J., Cook, B., Ishtiaq, S.: SLayer: Memory safety for systems-level code. In: Proc. CAV ’11

4. Brockschmidt, M., Ströder, T., Otto, C., Giesl, J.: Automated detection of non-termination and `NullPointerException` for JBC. In: Proc. FoVeOOS '11
5. Brockschmidt, M., Musiol, R., Otto, C., Giesl, J.: Automated termination proofs for Java programs with cyclic data. In: Proc. CAV '12
6. Brockschmidt, M., Cook, B., Fuhs, C.: Better termination proving through cooperation. In: Proc. CAV '13
7. Brotherston, J., Gorogiannis, N., Petersen, R.L.: A generic cyclic theorem prover. In: Proc. APLAS '12
8. Clang compiler: <http://clang.llvm.org>
9. Cousot, P., Cousot, R.: Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In: POPL '77
10. Dudka, K., Müller, P., Peringer, P., Vojnar, T.: Predator: A shape analyzer based on symbolic memory graphs (competition contribution). In: Proc. TACAS '14
11. Dutertre, B., de Moura, L.M.: The Yices SMT solver (2006), tool paper at <http://yices.csl.sri.com/tool-paper.pdf>
12. Falke, S., Kapur, D., Sinz, C.: Termination analysis of C programs using compiler intermediate languages. In: Proc. RTA '11
13. Falke, S., Merz, F., Sinz, C.: LLBMC: Improved bounded model checking of C using LLVM (competition contribution). In: Proc. TACAS '13
14. Harris, W.R., Lal, A., Nori, A.V., Rajamani, S.K.: Alternation for termination. In: Proc. SAS '10
15. Heizmann, M., Hoenicke, J., Leike, J., Podelski, A.: Linear ranking for linear lasso programs. In: Proc. ATVA '13
16. Kroening, D., Sharygina, N., Tsitovich, A., Wintersteiger, C.: Termination analysis with compositional transition invariants. In: Proc. CAV '10
17. Lattner, C., Adve, V.S.: LLVM: A compilation framework for lifelong program analysis & transformation. In: Proc. CGO '04
18. Löwe, S., Mandrykin, M., Wendler, P.: CPAchecker with sequential combination of explicit-value analyses and predicate analyses (comp. contr.). In: Proc. TACAS '14
19. Magill, S., Tsai, M.H., Lee, P., Tsay, Y.K.: Automatic numeric abstractions for heap-manipulating programs. In: Proc. POPL '10
20. de Moura, L.M., Bjørner, N.: Z3: An efficient SMT solver. In: Proc. TACAS '08
21. Moy, Y., Marché, C.: Modular inference of subprogram contracts for safety checking. *J. Symb. Comput.* 45(11), 1184–1211 (2010)
22. O'Hearn, P., Reynolds, J., Yang, H.: Local reasoning about programs that alter data structures. In: Proc. CSL '01
23. <http://fxr.watson.org/fxr/source/lib/libsa/strlen.c?v=OPENBSD>
24. Podelski, A., Rybalchenko, A.: ARMC: The logical choice for software model checking with abstraction refinement. In: Proc. PADL '07
25. Spoto, F., Mesnard, F., Payet, É.: A termination analyser for Java Bytecode based on path-length. *ACM TOPLAS* 32(3) (2010)
26. Ströder, T., Giesl, J., Brockschmidt, M., Frohn, F., Fuhs, C., Hensel, J., Schneider-Kamp, P.: Automated termination analysis for programs with pointer arithmetic. Tech. Rep. AIB 2014-05, available from [2] and [aib.informatik.rwth-aachen.de](http://aib.informatik.rwth-aachen.de)
27. *SV-COMP at TACAS 2014*: <http://sv-comp.sosy-lab.org/2014/>
28. Tsitovich, A., Sharygina, N., Wintersteiger, C.M., Kroening, D.: Loop summarization and termination analysis. In: Proc. TACAS '11
29. Urban, C.: The abstract domain of segmented ranking functions. In: Proc. SAS '13
30. Wikibooks C Programming: [http://en.wikibooks.org/wiki/C\\_Programming/](http://en.wikibooks.org/wiki/C_Programming/)
31. Zhao, J., Nagarakatte, S., Martin, M.M.K., Zdancewic, S.: Formalizing the LLVM IR for verified program transformations. In: Proc. POPL '12