

Verifying Pointer Safety for Programs with Unknown Calls

Chenguang Luo, Florin Craciun, Shengchao Qin, Guanhua He

Durham University, Durham DH1 3LE, UK

Wei-Ngan Chin

National University of Singapore

Abstract

We study the automated verification of pointer safety for heap-manipulating imperative programs with unknown procedure calls. Given a Hoare-style partial correctness specification $S = \{\text{Pre}\} C \{\text{Post}\}$ in separation logic, where the program C contains calls to some unknown procedure U , we infer a specification S_U for the unknown procedure U from the calling contexts. We show that the problem of verifying the program C against the specification S can be safely reduced to the problem of proving that the procedure U (once its code is available) meets the derived specification S_U . The expected specification S_U for the unknown procedure U is automatically calculated using an abduction-based shape analysis. We have also implemented a prototype system to validate the viability of our approach. Preliminary results show that the specifications derived by our tool fully capture the behaviors of the unknown code in many cases.

Key words: Inference, verification, abduction, separation logic

1. Introduction

While automated verification of memory safety remains a big challenge [14, 21], especially for substantial system software such as the Linux distribution and device drivers, significant advances have been seen recently on the automated verification of pointer safety for such software [6, 22]. The abduction-based compositional shape analysis [6] is able to calculate the pre/post specifications for procedures in a bottom-up approach, based on the call-dependency graph. It can verify the pointer safety of a large portion

Email addresses: {chenguang.luo, florin.craciun, shengchao.qin, guanhua.he}@durham.ac.uk (Chenguang Luo, Florin Craciun, Shengchao Qin, Guanhua He), chinwn@comp.nus.edu.sg (Wei-Ngan Chin).

of the Linux kernel and many device drivers manipulating shared mutable data structures. One issue that has not been dealt within their work is unknown procedure calls. In their SpaceInvader tool, unknown calls are currently ignored and replaced by the empty statement `skip` during the verification. This may lead to imprecise or unsound results in general. Our work here is to investigate this issue carefully and provide a better solution to the verification of pointer safety of programs with unknown calls.

Automated program verifiers usually require to have access to the entire given program, which in practice may not be completely available for various reasons. For instance, some programs (e.g. in C) may contain unknown calls that correspond to function pointers, some programs (e.g. in OO) may contain calls to interface methods whose actual implementations may not be available statically, or some programs may have calls to library procedures whose code is not available during verification. Other possible scenarios can be found in remote procedure calls such as COM/DCOM [20], mobile code and software upgrading, where program fragments may be unavailable at verification time. To deal with the verification of programs with unknown procedure calls, current automated program verifiers either

- ignore the unknown procedure calls, e.g. replacing them by `skip` [6], which can be unsound in general; or
- assume that the program and the unknown procedure have disjoint memory footprints so that the unknown call can be safely ignored due to the hypothetical frame rule [18]; however, this assumption does not hold in many cases; or
- use specification mining [1] to discover possible specifications for the (unknown part of the) program, which is performed dynamically by observing the execution traces and is not likely to be exhaustive for all possible program behaviors; or
- take into account all possible implementations for the unknown procedure [10, 12]. In general, there can be too many such candidates, making the verification almost impossible at compile time; or
- simply stop at the first unknown procedure call and provide an incomplete verification, which is obviously undesirable.

Approach and contributions. We propose a different approach in this paper to the verification of programs with unknown procedure calls. Given a specification $S = \{\text{Pre}\} C \{\text{Post}\}$ for the program C containing calls to an unknown procedure U , our solution is to proceed with the verification for the known fragments of C , and at the same time infer a specification S_U that is expected for the unknown procedure U based on the calling context(s). The problem of verifying the program C against the specification S can now be safely reduced to the problem of verifying the procedure U against the inferred specification S_U , provided that the verification of the known fragments does not cause any problems. The inferred specification is subject to a later verification when an implementation or a specification for the unknown procedure becomes available (e.g. at loading time in Java).

The intuition of our method to infer the unknown procedure’s specification can be divided into two steps. The first step is to analyze the code before the unknown procedure call to discover its precondition. The second is to analyze the code after the unknown call in order to discover its postcondition. For the second step, one might suggest to use a deductive backwards analysis, starting from the postcondition of the program being verified, to derive an expected postcondition for the unknown procedure. We cannot follow this suggestion due to a technical challenge: backwards reasoning over the separation

domain is simply too costly to implement ([6]). Therefore, we exploit an abductive forwards reasoning [6, 11] to derive the unknown procedure’s postcondition, and in this way the verification can be accomplished.

Our paper makes the following technical contributions:

- We propose a novel framework in separation logic for the verification of pointer safety for programs with unknown calls.
- The *top-down* feature of our approach can potentially benefit the general software development process. Given the specification for the caller procedure, it can be used to infer the specification for the callee procedures. This is a potentially beneficial complement for current *bottom-up* approaches ([13]).
- We have enhanced the abduction mechanism, resulting in an improved algorithm for the verification use. Compared with bi-abductive analysis [6] we have added additional reasoning rules (e.g. for the list tail matching) such that the use of inaccurate reasoning rules (like the “missing” rule) is avoided in many cases. We have also introduced a more practical partial order to judge the quality of different solutions for abduction.
- We have successfully incorporated the call-by-reference mechanism into the forward analyses in the separation domain.
- We have built a prototype system to test the viability and performance of our approach, and we have conducted some initial experimental studies to evaluate the precision of the results and the scalability of our system. Preliminary results show that our tool can derive expressive specifications which fully capture the behaviors of the unknown code in many cases.

Outline. Section 2 employs a motivating example to informally illustrate our main approach. Section 3 presents the programming language and the abstract domain for our analysis. Section 4 introduces our abductive reasoning. Section 5 defines two abstract semantics used in our verification. Section 6 depicts our verification algorithms. Experimental results are shown in Section 7, followed by some concluding remarks.

2. A Motivating Example

In this section we illustrate informally, via an example, how our analysis infers the specification for an unknown procedure. Our analysis makes use of a separation domain similar to the one used in the SpaceInvader tool [6, 9]. To keep the presentation simple, we use a small imperative language with both call-by-value and call-by-reference parameters for procedures. Formal details about the abstract domain and the language will be given in Section 3.

Example 1. Our goal is to verify the procedure `findLast` against the given pre/post specifications shown in Figure 1. The data structure `node { int data; node next }` defines a node in a linked list. The predicate $ls(x, y)$ used in the pre/post specifications as well as other places denotes a (possibly empty) list segment referred to by x and ended with the pointer y (i.e., y denotes the `next` field of the last node). Its formal definition is given in page 7.

According to the given specification, the procedure `findLast` takes in a non-empty linked list x and stores a reference to the last node of the list in the call-by-reference parameter z . Here we group call-by-value and call-by-reference parameters together and use semicolon `;` to separate them. Note that `findLast` calls an unknown procedure `unkProc` at line 4. □

```

// Given Specification:
// PrefindLast := ls(x, null) ∧ x≠null
// PostfindLast := ls(x, z) * z↦null
void findLast(node x; ref node z) {
0  node w, y;
0a // Δ := PrefindLast                Δ ⊢ x≠null
1  w := x.next;
1a // Δ := x↦w * ls(w, null) ∧ x≠null
2  if (w == null) z := x;
2a // Δ := x↦w ∧ x≠null ∧ w=null ∧ z=x
2b // ∃w, y · Δ ⊢ PostfindLast
3  else {
3a // Δ := x↦w * ls(w, null) ∧ x≠null ∧ w≠null
3b // H := Local(Δ, {x, y}) := x↦w * ls(w, null) ∧ x≠null ∧ w≠null
3c // R0 := Frame(Δ, {x, y}) := emp ∧ x≠null ∧ w≠null
4  unkProc(x; y);
4a // Δ := R0 * M0    M0 := (emp ∧ x=a ∧ y=b)    M := M0
4b // Δ ⊈ [y/x] PrefindLast
4c // Δ * [M1] ▷ [y/x] PrefindLast (s.t. Δ * M1 ⊢ [y/x] PrefindLast * true)
4d // M1 := ls(y, null) ∧ y≠null    M := M * M1
4e // Δ * M1 ⊢ [y/x] PrefindLast * R1    R1 := emp * x=a ∧ y=b
5  findLast(y; z);
5a // Δ := ([y/x] PostfindLast) * R1
5b // Δ ⊈ PostfindLast
5c // Δ * [M2] ▷ PostfindLast    M2 := ls(x, y)    M := M * M2
6  } }
6a // PreunkProc := ∃w · [a/x, b/y] H
6b // PostunkProc := [a/x, b/y] M

```

Fig. 1. Verification of `findLast` calling an unknown procedure `unkProc`.

For this example, the unknown call to `unkProc` performs essential modification to a local variable, so it can neither be regarded as `skip` nor dealt with using the hypothetical frame rule. The candidate implementations of `unkProc` are not available, preventing a verification against all possible candidates. However our approach is still applicable here as described below.

We conduct a symbolic execution ([4, 16]) on the procedure body starting with the precondition $\text{Pre}_{\text{findLast}}$ (line 0a). The results of our analysis (e.g. the abstract states) are marked as comments in the code. The analysis carries on as a standard forward shape analysis until the unknown procedure call at line 4.

At line 3, the current symbolic heap Δ is split into two disjoint parts: the local part H (line 3b) that is dependent on, and possibly mutated by, the unknown procedure; and the frame part R_0 (line 3c) that is not accessed by the unknown procedure. Intuitively, the local part of a symbolic heap w.r.t. a set of variables X is the part of the heap

reachable from variables in X (together with the aliasing information); while the frame part denotes the unreachable heap part (together with the aliasing information). For example, for a symbolic heap $\text{ls}(x, w) * \text{ls}(y, z) * \text{ls}(z, \text{null}) \wedge w=z$, its local part w.r.t. $\{x\}$ is $\text{ls}(x, w) * \text{ls}(z, \text{null}) \wedge w=z$, and its frame part w.r.t. $\{x\}$ is $\text{ls}(y, z) \wedge w=z$. We will have their formal definitions in Section 6.

We take H (line 3b) as a crude precondition for the unknown procedure, since it denotes the symbolic heap that is accessible, and hence potentially usable, by the unknown call. The frame part R_0 is not touched by the unknown call and will remain in the post-state, as shown in line 4a.

At line 4a, the abstract state after the unknown call consists of two parts: one is the aforesaid frame R_0 not accessed by the call, and the other is due to the procedure's postcondition which is unfortunately not available. Our next step is to discover the postcondition by examining (the requirements of) the code fragment after the unknown call by doing abductive reasoning (line 4a to 5c).

Initially, we assume the unknown procedure having an empty heap M_0 as its postcondition¹, and gradually discover the missing parts of the postcondition during the symbolic execution of the code fragment after the unknown call. To do that, our analysis keeps track of a pair (Δ, M) at each program point, where Δ refers to the current heap state, and M denotes the expected postcondition discovered so far for the unknown procedure. The notations M_i are used to represent parts of the discovered postcondition.

At line 5, the procedure `findLast` is called recursively. Since the current heap state does not satisfy the precondition of `findLast` (as shown in line 4b), the verification fails. However, this is not necessarily due to a program error; it may be due to the fact that the unknown call's postcondition is still unknown. Therefore, our analysis performs an abductive reasoning (line 4c) to infer the missing part M_1 for Δ such that $\Delta * M_1$ entails the precondition of `findLast` w.r.t. some substitution $[y/x]$. As shown in line 4d, M_1 is inferred to be $\text{ls}(y, \text{null}) \wedge y \neq \text{null}$, which is accumulated into M as part of the expected postcondition of the unknown procedure. (We will explain the details for abductive reasoning in Sec 4.) Now the heap state combined with the inferred M_1 meets the precondition of the procedure `findLast`, and also generates a residual frame heap R_1 (line 4e).

The heap state Δ immediately after the recursive call (line 5a) is formed by `findLast`'s postcondition and the frame R_1 , and it is expected to establish the postcondition of `findLast` for the overall verification to succeed. However, it does not (as shown in line 5b). Again this might be due to the fact that part of the unknown call's postcondition is still missing. Therefore, we perform another abductive reasoning (line 5c) to infer the missing M_2 as follows:

$$\text{ls}(y, z) * z \mapsto \text{null} \wedge x=a \wedge y=b * [M_2] \triangleright \text{ls}(x, z) * z \mapsto \text{null}$$

such that $\Delta * M_2$ entails $\text{Post}_{\text{findLast}}$. In this case, our abductor returns $M_2 := \text{ls}(x, y)$ as the result which is then added into M by separation conjunction, as shown in line 5c.

Finally, we generate the expected pre/post-specification for the unknown procedure (lines 6a and 6b). The precondition is obtained from the local pre-state of the unknown call, H , by replacing all variables that are aliases of a (or b) with the formal parameter a (or b). The postcondition is obtained from the accumulated abduction result, M ,

¹ Note that we introduce fresh logical variables a and b to record the value of x and y when `unkProc` returns.

after performing the same substitution. Our discovered specification for the unknown procedure $\text{unkProc}(a; b)$ is:

$$\begin{aligned} \text{Pre}_{\text{unkProc}} &:= \exists w \cdot a \mapsto w * \text{ls}(w, \text{null}) \wedge a \neq \text{null} \wedge w \neq \text{null} \\ \text{Post}_{\text{unkProc}} &:= \text{ls}(a, b) * \text{ls}(b, \text{null}) \wedge b \neq \text{null} \end{aligned}$$

The entire program is correct if unkProc meets the derived specification.

3. Programming Language and Abstract Domain

In this section, we first depict the syntax of a language in which programs may invoke unknown procedures, and then present the abstract domain for our analysis.

To focus only on key issues, we use a simple imperative language:

$$\begin{aligned} E &=_{df} x \mid \text{null} \\ b &=_{df} E_1 = E_2 \mid E_1 \neq E_2 \\ A[E] &=_{df} [E] := E_1 \mid \text{dispose}(E) \mid x := [E] \\ A &=_{df} x := E \mid x := \text{new}(E) \mid \text{skip} \\ C &=_{df} A[E] \mid A \mid f(\mathbf{x}; \mathbf{y}) \mid C_1; C_2 \mid \\ &\quad \text{if } b \text{ then } C_1 \text{ else } C_2 \text{ fi} \mid \text{while } b \text{ do } C \text{ od} \\ U &=_{df} \text{unkFn}(\mathbf{x}; \mathbf{y}) \mid \{ \text{unkFn}(\mathbf{x}_0; \mathbf{y}_0); C_1; \\ &\quad \text{unkFn}(\mathbf{x}_1; \mathbf{y}_1); C_2; \dots; C_{n-1}; \text{unkFn}(\mathbf{x}_n; \mathbf{y}_n) \} \mid \\ &\quad \text{if } b \text{ then } V \text{ else } C \text{ fi} \mid \text{if } b \text{ then } C \text{ else } V \text{ fi} \mid \\ &\quad \text{if } b \text{ then } V_1 \text{ else } V_2 \text{ fi} \mid \text{while } b \text{ do } V \text{ od} \\ V &=_{df} \{ C_1; U_{(\mathbf{x}_1; \mathbf{y}_1)}; C_2 \}_{(\mathbf{x}_0; \mathbf{y}_0)} \\ P &=_{df} \cdot \mid P; f(\mathbf{x}; \text{ref } \mathbf{y}) \{ \text{local } \mathbf{z}; C \} \mid \\ &\quad P; f(\mathbf{x}; \text{ref } \mathbf{y}) \{ \text{local } \mathbf{z}; V \} \end{aligned}$$

Note that expressions (E) are program variables to record memory locations in the heap, and all program variables are assumed of the same type, reference. The language has both heap sensitive ($A[E]$) and heap insensitive (A) atomic commands. The former requires access to the heap location referred to by E , while the latter does not. The command C contains calls to known procedures only, while the commands U and V comprise unknown procedure calls. Note also that our language allows both call-by-value and call-by-reference parameters for procedures, and for convenience, we group call-by-value parameters on the left and call-by-reference ones on the right, separated by a semicolon.

The unknown commands U and V specify the possible scenarios in which an unknown call ($\text{unkFn}(\mathbf{x}; \mathbf{y})$) may occur. Note that U and V may be annotated with two sets of variables, e.g. $(\mathbf{x}_0, \mathbf{y}_0)$ in the definition of V , where \mathbf{x}_0 denotes variables that can be accessed, but cannot be modified by V , and where \mathbf{y}_0 denotes variables that may be mutated by V . The same annotation applies to U . These annotations can be obtained automatically via a pre-processing phase of the analysis, by recording all the variables appearing in V as the right(left)-value of an assignment, and/or as a call-by-value(reference) parameter of a procedure invocation [17].

Example 2 (Unknown procedure and unknown block). The parse tree of `findLast`'s body (omitting local variable definition) in our motivating example is as in Figure 2. \square

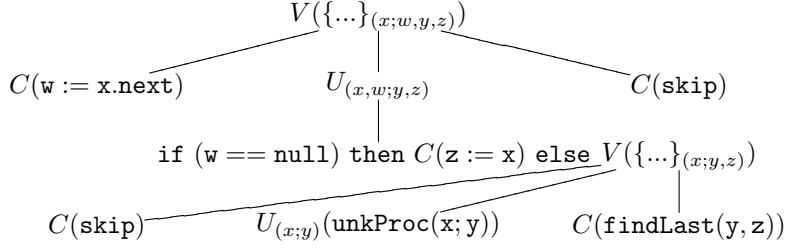


Fig. 2. The parse tree of `findLast`'s body.

A program P is composed of several procedures (with one of them as the entry method). A procedure can be totally known to the verifier (so that its body is composed by C), or it can contain an unknown block V .

We use a *memory model* similar to the standard one for separation logic [19]:

$$\begin{aligned}
 \text{Stack} &=_{df} (\text{Var} \cup \text{LVar} \cup \text{SVar}) \rightarrow \text{Val} \\
 \text{Val} &=_{df} \text{Loc} \cup \{\text{null}\} \\
 \text{Heap} &=_{df} \text{Loc} \rightarrow \text{Val} \\
 \text{State} &=_{df} \text{Stack} \times \text{Heap}
 \end{aligned}$$

The slight difference is that we have three (disjoint) set of variables: a finite set of program variables $\text{Var} = \{x, y, ..\}$, logical variables $\text{LVar} = \{x', y', ..\}$, and the special (logical) variables $\text{SVar} = \{a, b, ..\}$, with the last set reserved for unknown procedures for specification purpose. As usual, Loc is a countably infinite set of locations, subsumed by the set of values Val . The function Heap denotes a partial mapping from locations to values and a program state is a pair of stack and heap.

An abstract (program) state in our analysis is a *symbolic heap* representing a set of concrete heaps. It is defined as follows:

E	$=_{df} x \mid x' \mid a$	Expressions
Π	$=_{df} E_1 = E_2 \mid E_1 \neq E_2 \mid \text{true} \mid \Pi_1 \wedge \Pi_2$	Pure formulae
$\text{B}(E_1, E_2)$	$=_{df} E_1 \mapsto E_2 \mid \text{ls}(E_1, E_2)$	Basic separation predicates
Σ	$=_{df} \text{B}(E_1, E_2) \mid \text{true} \mid \text{emp} \mid \Sigma_1 * \Sigma_2$	Separation formulae
Δ	$=_{df} \Pi \wedge \Sigma$	Quantifier-free symbolic heaps
H	$=_{df} \exists \mathbf{x}. \Delta$	Symbolic heaps

The expressions (x, x' and a) correspond to the three kinds of variables (program, logical and special ones). Pure formulae Π express the aliasing information among expressions. The basic separation predicate $\text{B}(E_1, E_2)$ denotes either a singleton heap or a list segment [6, 9, 16]. A (possibly empty) list segment is inductively defined as follows:

$$\text{ls}(E_1, E_2) =_{df} (\text{emp} \wedge E_1 = E_2) \vee (\exists E_3. E_1 \mapsto E_3 * \text{ls}(E_3, E_2))$$

Here the list segment predicate is regarded as a built-in predicate in our abstract domain, as our abstract semantics and abduction are fine tuned for it. According to *SpaceInvader* [6] we may extend it with other predicates. A symbolic heap H is composed of a

pure formula Π and a separation formula Σ , possibly with existential quantifications over them. The separation formula Σ is formed by the predicates $B(E_1, E_2)$, **true** (arbitrary heaps) and **emp** (an empty heap) via separation conjunction [19]. We use **SH** to denote the set of all symbolic heaps and we will use the two terms “symbolic heap” and “abstract state” interchangeably.

The semantics of a symbolic heap H [3, 4] is defined by the satisfaction relationship $(s, h) \models H$ between a pair of stack and heap states (s, h) and H , where $s \in \text{Stack}$ and $h \in \text{Heap}$:

$$\begin{aligned}
& \llbracket x \rrbracket s =_{df} s(x) & \llbracket \text{null} \rrbracket s =_{df} \text{null} \\
(s, h) \models E_1 = E_2 & =_{df} \llbracket E_1 \rrbracket s = \llbracket E_2 \rrbracket s & (s, h) \models E_1 \neq E_2 & =_{df} \llbracket E_1 \rrbracket s \neq \llbracket E_2 \rrbracket s \\
(s, h) \models \Pi_1 \wedge \Pi_2 & =_{df} (s, h) \models \Pi_1 \text{ and } (s, h) \models \Pi_2 \\
(s, h) \models \text{true} & =_{df} \text{always} \\
(s, h) \models \text{emp} & =_{df} h = \emptyset & (s, h) \models E_1 \mapsto E_2 & =_{df} h = \{ \llbracket E_1 \rrbracket s \mapsto \llbracket E_2 \rrbracket s \} \\
(s, h) \models \Sigma_1 * \Sigma_2 & =_{df} h = h_1 \cup h_2 \text{ and } (s, h_1) \models \Sigma_1 \text{ and } (s, h_2) \models \Sigma_2 \\
& \text{where } h_1 \text{ and } h_2 \text{ are domain-disjoint}
\end{aligned}$$

and the semantics of $\text{ls}(E_1, E_2)$ can be obtained from its inductive definition over the abstract domain. Based on such semantics, the entailment relationship between two symbolic heaps H_1 and H_2 is straightforward:

$$H_1 \vdash H_2 =_{df} \forall s, h. (s, h) \models H_1 \text{ implies } (s, h) \models H_2$$

which is used to check whether one symbolic heap is stronger than another in our analysis.

4. Abduction

As shown in Example 1 (Section 2), when analyzing the code after an unknown call, due to the lack of information about the unknown procedure, it is possible that the current state is too weak to meet the required precondition for the next instruction. As a consequence, the symbolic execution fails. A technique called abduction (or abductive reasoning) [6, 11] can be used to discover a symbolic heap M to make the entailment $\Delta * M \vdash H * R$ succeed, when the entailment $\Delta \vdash H * R$ fails. Here R denotes the (automatically computed) frame part. For instance, the entailments at line 4b and line 5b failed in Example 1, and in both cases, the abduction algorithm was called to find the missing M .

One problem in abduction is that there can be many solutions of M for the entailment $\Delta * M \vdash H * R$ to succeed. For instance, **false** can be a solution but should be avoided where possible. As another example, for the entailment $\text{ls}(y, z) * M \vdash \text{ls}(x, z) * R$ (a similar one was at line 5c in Example 1 in Sec 2), an abductive reasoning may return, for example, two different solutions: $M = \text{ls}(x, y)$ with $R = \text{emp}$ or $M = \text{ls}(x, z)$ with $R = \text{ls}(y, z)$. A partial order \preceq over symbolic heaps was given by Calcagno et al. [6] to make the selection from different solutions:

$$\begin{aligned}
M \preceq M' & =_{df} (M' \vdash M * \text{true} \wedge M \not\vdash M' * \text{true}) \vee \\
& (M' \vdash M * \text{true} \wedge M \vdash M' * \text{true} \wedge M' \vdash M)
\end{aligned}$$

Intuitively speaking, consider two symbolic heaps M and M' , if M' entails M (possibly extended by some frame), then M is less than M' under \preceq . For example, we have $\text{ls}(x, y) \preceq \text{ls}(x, y) * \text{ls}(y, z)$, and $\text{emp} \vee \exists x, y. \text{ls}(x, y) \preceq \text{emp}$.

With such partial order \preceq , we can judge the quality of different solutions for abduction by always choosing the least one according to \preceq .

Two solutions can not be distinguished when they are not comparable under this partial order. However, we may still prefer one solution over another in this case, as will be discussed in Section 6. We would expect that the solution to incur as few free variables as possible in the frame part (R). For this reason, we define a new order as follows:

$$\begin{aligned} M \preceq_H^{\Delta} M' =_{df} & M \preceq M' \vee (M \not\preceq M' \wedge M' \not\preceq M \wedge \\ & \Delta * M \vdash H * R \wedge \Delta * M' \vdash H * R' \wedge \\ & |\text{fv}(R)| \leq |\text{fv}(R')|) \end{aligned}$$

Note that our partial order is defined over the set of solutions for one abduction. Intuitively, given two solutions which are not comparable w.r.t. \preceq , our partial order may still be able to compare them, if one incurs fewer free variables in the frame R than the other. To ensure that the abductor will choose the former, but not the latter, we enrich the abductor [6] with some new rules (right column below) so that it attempts to introduce less free variables in the frame, where possible:

$$\begin{array}{l} \frac{\Delta * [M] \triangleright \Delta' \quad \Delta * B(E, E') \not\vdash \text{false}}{\Delta * [M * B(E, E')] \triangleright \Delta' * B(E, E')} \text{missing} \quad \frac{(E_0 = E_1 \wedge \Delta) * [M] \triangleright \exists \mathbf{y}'. \Delta'}{\Delta * E_0 \mapsto E * [\exists \mathbf{x}'. E_0 = E_1 \wedge M] \triangleright \exists \mathbf{x}' \mathbf{y}'. \Delta' * E_1 \mapsto E} \text{t-match} \\ \\ \frac{(E \neq E_0 \wedge \Delta * \text{ls}(x, E_0)) * [M] \triangleright \Delta'}{\Delta * \text{ls}(E, E_0) * [\exists \mathbf{x}. \Delta' * E \mapsto x]} \text{h-left} \quad \frac{E \neq E_0 \wedge \Delta * \text{ls}(E_0, x') * [M] \triangleright \exists \mathbf{y}'. \Delta'}{\Delta * \text{ls}(E_0, E) * [\exists \mathbf{x}' \mathbf{y}'. \Delta' * x' \mapsto E]} \text{t-left} \\ \\ \frac{\Delta * [M] \triangleright \Delta' * \text{ls}(E_0, E_1)}{\Delta * B(E, E_0) * [M] \triangleright \Delta' * \text{ls}(E, E_1)} \text{h-right} \quad \frac{\Delta * [M] \triangleright \exists \mathbf{x}'. \Delta' * \text{ls}(E_0, E_1)}{\Delta * B(E_0, E) * [M] \triangleright \exists \mathbf{x}'. \Delta' * B(E_1, E)} \text{t-right} \end{array}$$

where the rules from Calcagno et al. [6] are on the left and our new ones are on the right. Compared with their work, when the matching from the head of a list segment fails, our abductor also tries to match from the tail, instead of directly applying **missing** rule to introduce a new separation predicate, which could bring in more free variables in the frame. In our case, the rule **t-match** introduces aliasing information between two heads of pointing-to relationships in the premise and conclusion of the abduction, respectively. The last two rules try to match from the end of a basic separation predicate; if the matching succeeds then the matched part will be dropped from both sides to reduce the complexity of the predicates. For instance, in Example 1 in Section 2, for the last abduction

$$\exists y. \text{ls}(y, z) * z \mapsto \text{null} * [M] \triangleright \text{ls}(x, z) * z \mapsto \text{null}$$

their abductor will return $\text{ls}(x, z)$ as M , by trying to match the heads of the list segments on both sides of \triangleright , and adding the missing heap part to the left hand side once the matching fails. This solution is not optimal in the light that it introduces a free variable z in the frame. Comparatively, our abduction tries to match also from the tail of a list segment, finding $\text{ls}(x, y)$ as the result for the aforesaid abduction, which is a minimal solution under our partial order \preceq_H^{Δ} .

We prefer to find solutions that are (potentially locally) minimal with respect to \preceq_H^Δ and consistent. However, such solutions are generally not easy to compute and can incur excess cost (with additional disjunction in the analysis). Therefore, our abductive inference is designed more from a practical perspective to discover anti-frames that should be suitable as specifications for unknown procedures, and the order \preceq_H^Δ is to state the decision choices of our implementation of the abduction. This is proven in Sec 7 with the experiments we conduct.

5. Abstract Semantics

This section introduces two kinds of abstract semantics we use to analyze the program: an underlying semantics from the local shape analysis ([9]) and another semantics based on both the first one and abduction from the composite shape analysis ([6]).

We denote the specifications of a procedure $f(\mathbf{x}; \mathbf{y})$ as a subset of $\text{SH} \times (\text{SH} \cup \{\top\})$ (where \top stands for a fault abstract state). Note that for a call-by-reference parameter x , both $\text{old}(x)$ and x may occur in a postcondition with the former referring to the value of x in the pre-state (as in JML [15]). To illustrate, a postcondition $\text{ls}(\text{old}(x), x)$ means that there is a list segment beginning with the initial value of x (present in the precondition) to the final x when the procedure returns. The set of all specifications for all procedures is defined as

$$\text{AllSpec} =_{df} \mathcal{P}(\text{Name} \times \mathcal{P}(\text{SH} \times (\text{SH} \cup \{\top\})))$$

where Name refers to the set of all function names. Then our underlying abstract semantics' type is defined as

$$\llbracket C \rrbracket : \text{AllSpec} \rightarrow \mathcal{P}^\top(\text{SH}) \rightarrow \mathcal{P}^\top(\text{SH})$$

where $\mathcal{P}^\top(\text{SH})$ stands for $\mathcal{P}(\text{SH} \cup \{\top\})$. Given a program C , a specification table $\mathcal{T} \in \text{AllSpec}$, a set of abstract states S , $\llbracket C \rrbracket_{\mathcal{T}} S$ returns another set of abstract states.

Example 3 (Underlying semantics). For the `findLast` in our motivating example, suppose the specification table \mathcal{T} is

$$\{(\text{findLast}(\mathbf{a}, \mathbf{b}), \{(\text{ls}(a, \text{null}) \wedge a \neq \text{null}, \text{ls}(a, b) * b \mapsto \text{null})\})\}$$

Then we know the symbolic execution $\llbracket \text{findLast}(x, y) \rrbracket_{\mathcal{T}} \{\text{ls}(x, \text{null}) \wedge x \neq \text{null}\}$ will give $\{\text{ls}(x, y) * y \mapsto \text{null}\}$ as result. \square

The basic transition functions below form the foundation of the first underlying semantics. With one symbolic heap as input, they return either another symbolic heap or a set of symbolic heaps:

$$\begin{aligned} \text{rearr}(E) &=_{df} \text{SH} \rightarrow \mathcal{P}^\top(\text{SH}[E]) && \text{Rearrangement} \\ \text{exec}(A[E]) &=_{df} \text{SH}[E] \rightarrow \text{SH} \cup \{\top\} && \text{Heap-sensitive execution} \\ \text{exec}(A) &=_{df} \text{SH} \rightarrow \text{SH} && \text{Heap-insensitive execution} \\ \text{abs} &=_{df} \text{SH} \rightarrow \text{SH} && \text{Abstraction} \end{aligned}$$

where $\text{SH}[E]$ denotes a set of symbolic heaps in which each element has E exposed as the head in one of its pointing-to conjuncts ($E \mapsto F$). Here $\text{rearr}(E)$ attempts to unroll

the shape beginning with E to expose it as the head of a pointing-to predicate, say, E points to another expression. (If such unrolling fails then it will return $\{\top\}$.) This is a preparation for the second transition function $\text{exec}(A[E])$, as it tries to perform some dereference of E . The third function, like the semantics for stack-based variable assignment and heap allocation, does not require such exposure of E , and thus is called heap-insensitive compared with the second one. The last transition function conducts a rolling over the shapes, to eliminate unimportant cutpoints to ensure termination for our verification.

The rules for the basic transition functions are adopted from Distefano et al. [9], where the logical variable x'' is always fresh.

Below are the rules for rearrangement, where we try to find an explicit pointing-to beginning with the expression E to be unrolled (or its alias). If such unrolling fails then $\{\top\}$ is returned.

$$\begin{aligned}
\text{rearr}(E) (\exists \mathbf{x}' . \Pi \wedge \Delta * \text{ls}(G, F)) &=_{df} \mathbf{if} \ \Pi \vdash G=E \ \mathbf{then} \\
&\quad \{\exists \mathbf{x}' . \Pi \wedge E=F \wedge \Delta, \exists \mathbf{x}' . \Pi \wedge \Delta * E \mapsto F, \\
&\quad \exists \mathbf{x}', x'' . \Pi \wedge \Delta * E \mapsto x'' * \text{ls}(x'', F)\} \ \mathbf{else} \ \{\top\} \\
| (\exists \mathbf{x}' . \Pi \wedge \Delta * G \mapsto F) &=_{df} \mathbf{if} \ \Pi \vdash G=E \ \mathbf{then} \\
&\quad \{\exists \mathbf{x}' . \Pi \wedge \Delta * E \mapsto F\} \ \mathbf{else} \ \{\top\} \\
| (\exists \mathbf{x}' . \Pi \wedge \Delta) &=_{df} \{\top\}
\end{aligned}$$

Below are rules for symbolic execution to reflect the effects, over heap or not, of atomic commands.

$$\begin{aligned}
\text{exec} (x := E) (\exists \mathbf{x}' . \Pi \wedge \Delta) &=_{df} \exists \mathbf{x}' . x=[x''/x]E \wedge [x''/x](\Pi \wedge \Delta) \\
\text{exec} (x := [E]) (\exists \mathbf{x}' . \Pi \wedge \Delta * E \mapsto F) &=_{df} \exists \mathbf{x}' . x=[x''/x]F \wedge [x''/x](\Pi \wedge \Delta * E \mapsto F) \\
| (\exists \mathbf{x}' . \Pi \wedge \Delta) &=_{df} \top \\
\text{exec} ([E] := G) (\exists \mathbf{x}' . \Pi \wedge \Delta * E \mapsto F) &=_{df} \exists \mathbf{x}' . \Pi \wedge \Delta * E \mapsto G \\
| (\exists \mathbf{x}' . \Pi \wedge \Delta) &=_{df} \top \\
\text{exec} (x := \text{new}(E)) (\exists \mathbf{x}' . \Pi \wedge \Delta) &=_{df} \exists \mathbf{x}' . [x''/x](\Pi \wedge \Delta) * x \mapsto [x''/x]E \\
\text{exec} (\text{dispose}(E)) (\exists \mathbf{x}' . \Pi \wedge \Delta * E \mapsto F) &=_{df} \exists \mathbf{x}' . \Pi \wedge \Delta \\
| (\exists \mathbf{x}' . \Pi \wedge \Delta) &=_{df} \top
\end{aligned}$$

Below are definitions of rules for abstraction. Their rationale is to remove any heap garbage from the current state and eliminate logical cutpoints that are neither shared nor essential in denoting a cyclic list over the heap. In our semantics, when the abstraction is actually performed, the six rules are applied several rounds, with each round from the top one to the bottom, until the abstracted state does not change any more. Its termination

is discussed later in Section 6.2.

$$\begin{array}{c}
\frac{}{\exists \mathbf{z}'. E=x' \wedge \Pi \wedge \Sigma \rightsquigarrow \exists \mathbf{z}'. [E/x']\Pi \wedge \Sigma} \text{Equality} \\
\frac{x' \notin \text{LVar}(\Sigma)}{\exists \mathbf{z}'. E \neq x' \wedge \Pi \wedge \Sigma \rightsquigarrow \exists \mathbf{z}'. \Pi \wedge \Sigma} \text{Disequality} \\
\frac{x' \notin \text{LVar}(\Pi, \Sigma)}{\exists \mathbf{z}'. \Pi \wedge \Sigma * \mathbf{B}(x', E) \rightsquigarrow \exists \mathbf{z}'. \Pi \wedge \Sigma * \mathbf{true}} \text{Junk1} \\
\frac{x', y' \notin \text{LVar}(\Pi, \Sigma)}{\exists \mathbf{z}'. \Pi \wedge \Sigma * \mathbf{B}(x', y') * \mathbf{B}(y', x') \rightsquigarrow \exists \mathbf{z}'. \Pi \wedge \Sigma * \mathbf{true}} \text{Junk2} \\
\frac{x' \notin \text{LVar}(\Pi, \Sigma, E, F) \quad \Pi \vdash F = \mathbf{null}}{\exists \mathbf{z}'. \Pi \wedge \Sigma * \mathbf{B}_1(E, x') * \mathbf{B}(x', F) \rightsquigarrow \exists \mathbf{z}'. \Pi \wedge \Sigma * \text{ls}(E, \mathbf{null})} \text{Abs1} \\
\frac{x' \notin \text{LVar}(\Pi, \Sigma, E, F, E_1, F_1) \quad \Pi \vdash F = E_1}{\exists \mathbf{z}'. \Pi \wedge \Sigma * \mathbf{B}_1(E, x') * \mathbf{B}(x', F) * \mathbf{B}(E_1, F_1) \rightsquigarrow \exists \mathbf{z}'. \Pi \wedge \Sigma * \text{ls}(E, F) * \mathbf{B}(E_1, F_1)} \text{Abs2}
\end{array}$$

A lifting function p^\dagger is defined over partial functions on symbolic heaps to lift their domains and ranges to a powerset of SH, plus \top :

$$p^\dagger S =_{df} \{\top \mid \top \in S\} \cup \{p H \mid H \in S \setminus \{\top\}\}$$

and this function is overloaded for `rearr` only to lift its domain to $\mathcal{P}^\top(\text{SH})$:

$$\text{rearr}(E)^\dagger S =_{df} \{\top \mid \top \in S\} \cup \left(\bigcup_{H \in S \setminus \{\top\}} \text{rearr}(E) H \right)$$

The following function, `filt`, is to filter out any symbolic heap that does not satisfy the given aliasing constraint:

$$\text{filt}(E=F)(H) =_{df} \text{if } H \not\vee E \neq F \text{ then } H \wedge E=F \text{ else undefined}$$

$$\text{filt}(E \neq F)(H) =_{df} \text{if } H \not\vee E=F \text{ then } H \wedge E \neq F \text{ else undefined}$$

and the program constructors' semantics are based on the atomic ones defined above:

$$\begin{array}{ll}
\llbracket \mathbf{b} \rrbracket_{\mathcal{T}} S & =_{df} \text{filt}(\mathbf{b})^\dagger S \\
\llbracket A[E] \rrbracket_{\mathcal{T}} S & =_{df} \text{abs}^\dagger \circ \text{exec}(A[E])^\dagger \circ \text{rearr}(E)^\dagger S \\
\llbracket A \rrbracket_{\mathcal{T}} S & =_{df} \text{abs}^\dagger \circ \text{exec}(A)^\dagger S \\
\llbracket C_1; C_2 \rrbracket_{\mathcal{T}} S & =_{df} \llbracket C_2 \rrbracket_{\mathcal{T}} \circ \llbracket C_1 \rrbracket_{\mathcal{T}} S \\
\llbracket \text{if } \mathbf{b} \text{ then } C_1 \text{ else } C_2 \text{ fi} \rrbracket_{\mathcal{T}} S & =_{df} \llbracket \mathbf{b}; C_1 \rrbracket_{\mathcal{T}} S \cup \llbracket \neg \mathbf{b}; C_2 \rrbracket_{\mathcal{T}} S \\
\llbracket \text{while } \mathbf{b} \text{ do } C \text{ od} \rrbracket_{\mathcal{T}} S & =_{df} \llbracket \neg \mathbf{b} \rrbracket_{\mathcal{T}} (\text{fix } \lambda S'. S \cup \llbracket \mathbf{b}; C \rrbracket_{\mathcal{T}} S')
\end{array}$$

Finally we need to adapt the underlying abstract semantics to procedure invocation:

$$\begin{aligned} \llbracket f(\mathbf{x}; \mathbf{y}) \rrbracket_{\mathcal{T}} S =_{df} & \{ [\mathbf{x}/\mathbf{a}, \mathbf{y}/\mathbf{b}, \mathbf{y}'/\text{old}(\mathbf{b})] \text{Post} * [\mathbf{y}'/\mathbf{y}] R \mid \Delta \in S \wedge \\ & (\text{Pre}, \text{Post}) \in \text{Spec}_f \wedge \Delta \vdash [\mathbf{x}/\mathbf{a}, \mathbf{y}/\mathbf{b}] \text{Pre} * R \} \\ & \text{where } (f, \text{Spec}_f) \in \mathcal{T} \text{ and } \mathbf{y}' \text{ are fresh} \end{aligned}$$

where \mathcal{T} is an element of AllSpec . We assume here the formal parameters in f 's specifications are expressed in terms of \mathbf{a} , \mathbf{b} and $\text{old}(\mathbf{b})$, respectively. We first check whether the procedure's precondition is satisfied by the current state (subject to substitution), and then replace it with the (substituted) postcondition to indicate the effect of the procedure call.

Next we define the abstract semantics with abduction used in our verification to discover the effect of unknown procedure calls.

$$\llbracket C \rrbracket^{\Delta} : \text{AllSpec} \rightarrow \mathcal{P}^{\top}(\text{SH} \times \text{SH}) \rightarrow \mathcal{P}^{\top}(\text{SH} \times \text{SH})$$

Here $\mathcal{P}^{\top}(\text{SH} \times \text{SH})$ denotes $\mathcal{P}((\text{SH} \cup \{\top\}) \times (\text{SH} \cup \{\top\}))$. Given a specification table $\mathcal{T} \in \text{AllSpec}$, each element of the input (or output) for $\llbracket C \rrbracket_{\mathcal{T}}^{\Delta}$ is a pair of two symbolic heaps, of which the first denotes the current program state, and the second stands for the abduction result. In our framework, this semantics is used to accumulate the discovered effect of unknown procedure calls into the second symbolic heap in the pair.

Example 4 (Abstract semantics with abduction). Consider the same setting for `findLast` in Example 3. For the semantics with abduction $\llbracket \text{findLast}(\mathbf{x}, \mathbf{y}) \rrbracket_{\mathcal{T}}^{\Delta} \{ \text{ls}(x, \text{null}) \}$, we will have $\{ (\text{ls}(x, y) * y \mapsto \text{null}, x \neq \text{null}) \}$ as result. Here in order to obtain the first component of the pair as the final state $(\text{ls}(x, y) * y \mapsto \text{null})$, the second component $(x \neq \text{null})$ must be added to its corresponding input state $\text{ls}(x, \text{null})$. \square

This semantics also consists of the basic transition functions which composes the atomic instructions' semantics and then the program constructors' semantics. Here the basic transition functions are lifted as follows:

$$\begin{aligned} \text{Rearr}(E)(H, M) &=_{df} \\ & \text{let } \mathcal{H} = \text{rearr}(E)(H) \text{ and } S = \{ (H', M) \mid H' \in \mathcal{H} \cap \text{SH} \} \\ & \text{in if } (\top \notin \mathcal{H}) \text{ then } S \\ & \quad \text{else if } (H \vdash E = a \text{ for some } a \in \text{SVar}) \text{ and } (M \not\vdash a \mapsto x' \text{ for fresh } x' \in \text{LVar}) \\ & \quad \text{then } S \cup \{ (H * E \mapsto x', M * a \mapsto x') \} \text{ else } S \cup \{ \top \} \\ \text{Exec}(A)(H, M) &=_{df} \text{let } \mathcal{H} = \text{exec}(A)(H) \text{ in } \{ (H', M) \mid H' \in \mathcal{H} \} \cup \{ \top \mid \top \in \mathcal{H} \} \\ & \quad \text{where } A \text{ is } [E] := G \text{ or } \text{dispose}(E) \\ \text{Exec}(x := E)(\exists \mathbf{x}'. \Pi \wedge \Delta, M) &=_{df} (\exists \mathbf{x}'. x = \sigma E \wedge \sigma(\Pi \wedge \Delta), \sigma M) \\ \text{Exec}(x := [E])(\exists \mathbf{x}'. \Pi \wedge \Delta * E \mapsto F) &=_{df} (\exists \mathbf{x}'. x = \sigma F \wedge \sigma(\Pi \wedge \Delta * E \mapsto F), \sigma M) \\ & \quad \mid (\exists \mathbf{x}'. \Pi \wedge \Delta) =_{df} \text{if } \exists \mathbf{x}'. \Pi \wedge \Delta * E \mapsto F \vdash \text{false} \text{ then } (\top, \sigma M) \\ & \quad \quad \text{else } (\exists \mathbf{x}'. x = \sigma F \wedge \sigma(\Pi \wedge \Delta * E \mapsto F), \sigma(M * E \mapsto F)) \\ \text{Exec}(x := \text{new}(E))(\exists \mathbf{x}'. \Pi \wedge \Delta) &=_{df} (\exists \mathbf{x}'. \sigma(\Pi \wedge \Delta) * x \mapsto \sigma E, \sigma M) \\ \text{Abs}(H, M) &=_{df} (\text{abs}(H), \text{abs}(M)) \end{aligned}$$

In the lifted `Rearr`, in case of a rearrangement failure, we utilize abduction to discover a pointing-to in the current heap, followed by another attempt of rearrangement. If it succeeds, then the abduction result is confirmed and added to the current state. Note that in the definition of `Exec` we need to treat variable assignments specifically. As can be seen in the rules, when x is assigned to a new value, the original value is still preserved in a fresh logical variable x'' with a substitution $\sigma = [x''/x]$. This allows us to keep the connection among the history values of a variable and its latest value, which may be essential as a link from the unknown procedure's postcondition to its caller's postcondition. This is because the unknown procedure's postcondition may refer to one of such variable's historical values, and its caller's postcondition can count on that variable's final value. In that case, the recorded values in the abduction results will serve as a connection among these abstract states. `Abs` simply performs abstraction over both H and M .

The filter function `Filt` now only works on the first symbolic heap of a pair:

$$\text{Filt}(\mathbf{b})(H, M) =_{df} (\text{filt}(\mathbf{b})(H), M)$$

and the abstract semantics for the program constructors is as follows:

$$\begin{aligned} \llbracket \mathbf{b} \rrbracket_{\mathcal{T}}^{\mathbf{A}} S &=_{df} \text{Filt}(\mathbf{b})^{\dagger} S \\ \llbracket A[E] \rrbracket_{\mathcal{T}}^{\mathbf{A}} S &=_{df} \text{Abs}^{\dagger} \circ \text{Exec}(A[E])^{\dagger} \circ \text{Rearr}(E)^{\dagger} S \\ \llbracket A \rrbracket_{\mathcal{T}}^{\mathbf{A}} S &=_{df} \text{Abs}^{\dagger} \circ \text{Exec}(A)^{\dagger} S \\ \llbracket C_1; C_2 \rrbracket_{\mathcal{T}}^{\mathbf{A}} S &=_{df} \llbracket C_2 \rrbracket_{\mathcal{T}}^{\mathbf{A}} \circ \llbracket C_1 \rrbracket_{\mathcal{T}}^{\mathbf{A}} S \\ \llbracket \text{if } \mathbf{b} \text{ then } C_1 \text{ else } C_2 \text{ fi} \rrbracket_{\mathcal{T}}^{\mathbf{A}} S &=_{df} \llbracket \mathbf{b}; C_1 \rrbracket_{\mathcal{T}}^{\mathbf{A}} S \cup \llbracket \neg \mathbf{b}; C_2 \rrbracket_{\mathcal{T}}^{\mathbf{A}} S \\ \llbracket \text{while } \mathbf{b} \text{ do } C \text{ od} \rrbracket_{\mathcal{T}}^{\mathbf{A}} S &=_{df} \llbracket \neg \mathbf{b} \rrbracket_{\mathcal{T}}^{\mathbf{A}} (\text{Ifix } \lambda S'. S \cup \llbracket \mathbf{b}; C \rrbracket_{\mathcal{T}}^{\mathbf{A}} S') \end{aligned}$$

Same as above, at last is the semantics for procedure invocation with abduction:

$$\begin{aligned} \llbracket f(\mathbf{x}; \mathbf{y}) \rrbracket_{\mathcal{T}}^{\mathbf{A}} S &=_{df} \{ ([\mathbf{x}/\mathbf{a}, \mathbf{y}/\mathbf{b}, \mathbf{y}'/\text{old}(\mathbf{b})] \text{Post} * [\mathbf{y}'/\mathbf{y}] R, [\mathbf{y}'/\mathbf{y}] M) \mid (\Delta, M) \in S \wedge \\ &\quad (\text{Pre}, \text{Post}) \in \text{Spec}_f \wedge \Delta \vdash [\mathbf{x}/\mathbf{a}, \mathbf{y}/\mathbf{b}] \text{Pre} * R \} \cup \\ &\quad \{ ([\mathbf{x}/\mathbf{a}, \mathbf{y}/\mathbf{b}, \mathbf{y}'/\text{old}(\mathbf{b})] \text{Post} * [\mathbf{y}'/\mathbf{y}] R, [\mathbf{y}'/\mathbf{y}] (M * M')) \mid \\ &\quad (\Delta, M) \in S \wedge (\text{Pre}, \text{Post}) \in \text{Spec}_f \wedge \Delta \not\vdash [\mathbf{x}/\mathbf{a}, \mathbf{y}/\mathbf{b}] \text{Pre} * R \wedge \\ &\quad \Delta * [M'] \triangleright [\mathbf{x}/\mathbf{a}, \mathbf{y}/\mathbf{b}] \text{Pre} \} \\ &\text{where } (f, \text{Spec}_f) \in \mathcal{T} \text{ and } \mathbf{y}' \text{ are fresh} \end{aligned}$$

where \mathcal{T} is an element of `AllSpec`. We also use \mathbf{a} , \mathbf{b} and `old`(\mathbf{b}) for formal parameters in the specification. This rule distinguishes two cases: 1. the current state in S is sufficiently strong to entail `Pre` subject to some substitutions, and corresponding postcondition is established; 2. the current state does not entail substituted `Pre`, and abduction is applied with further requirements on the current state accumulated (M'). It will combine the results from both cases to update the state after the procedure invocation and continue with the symbolic execution.

6. Verification

Based on the abstract semantics defined in the last section, we present in this section our algorithms for the verification of programs with unknown calls.

6.1. Main Verification Algorithm

1. Main algorithm. Our verification algorithm given in Figure 3 attempts to verify the body of the current procedure (the third input, comprising an unknown command U) against the given specifications (the second input). The first input gives a set of known procedure specifications, which are necessary as the current procedure may invoke known procedures apart from unknown ones. If the verification succeeds, it returns specifications that are expected for all unknown procedures invoked within U for the whole verification to succeed. If it fails, we know that the current procedure cannot meet one or more given specifications, no matter what specifications are given to the invoked unknown procedures. Returned specifications will be expressed using special variables \mathbf{a}, \mathbf{b} , etc. as in the earlier example.

For each specification (Pre, Post) to verify against (line 2), the algorithm works in three steps. Based on the underlying semantics mentioned earlier, it first computes the post-states of C_1 (i.e. S_0) from Pre (line 3), from which it extracts the preconditions for $U_{(\mathbf{x};\mathbf{y})}$ using the function Local . Intuitively, it extracts the part of each Δ_1 reachable from the variables that may be accessed by U , namely, \mathbf{x} and \mathbf{y} (line 6). Here $\text{fv}(\Delta)$ stands for all free (program and logical) variables occurring in Δ . The function $\text{Local}(\Pi \wedge \Sigma, \{\mathbf{x}\})$ is defined as follows:

$$\text{Local}(\Pi \wedge \Sigma, \{\mathbf{x}\}) =_{df} \exists \text{fv}(\Pi \wedge \Sigma) \setminus \text{ReachVar}(\Pi \wedge \Sigma, \{\mathbf{x}\}) \cdot \\ \Pi * \text{ReachHeap}(\Pi \wedge \Sigma, \{\mathbf{x}\})$$

where $\text{ReachVar}(\Pi \wedge \Sigma, \{\mathbf{x}\})$ is the minimal set of variables reachable from $\{\mathbf{x}\}$:

$$\{\mathbf{x}\} \cup \{z_2 \mid \exists z_1, \Pi_1 \cdot z_1 \in \text{ReachVar}(\Pi \wedge \Sigma, \{\mathbf{x}\}) \wedge \Pi = (z_1 = z_2) \wedge \Pi_1\} \cup \\ \{z_2 \mid \exists z_1, \Sigma_1 \cdot z_1 \in \text{ReachVar}(\Pi \wedge \Sigma, \{\mathbf{x}\}) \wedge \Sigma = \mathbf{B}(z_1, z_2) * \Sigma_1\} \subseteq \text{ReachVar}(\Pi \wedge \Sigma, \{\mathbf{x}\})$$

where $\mathbf{B}(z_1, z_2)$ stands for either $z_1 \mapsto z_2$ or $\text{ls}(z_1, z_2)$. And the formula $\text{ReachHeap}(\Pi \wedge \Sigma, \{\mathbf{x}\})$ denotes the part of Σ reachable from $\{\mathbf{x}\}$ and is formally defined as the *-conjunction of the following set of formulae:

$$\{\Sigma_1 \mid \exists z_1, z_2, \Sigma_2 \cdot z_1 \in \text{ReachVar}(\Pi \wedge \Sigma, \{\mathbf{x}\}) \wedge \Sigma = \Sigma_1 * \Sigma_2 \wedge \Sigma_1 = \mathbf{B}(z_1, z_2)\}$$

The second step is to symbolically execute C_2 , using the abstract semantics with abduction, to discover the postconditions for $U_{(\mathbf{x};\mathbf{y})}$ (lines 8–10). At line 8, since we know nothing about U , we take emp as the post-state of U . Therefore, the initial state for the symbolic execution of C_2 is simply the frame part of state not touched by U . Here Frame is formally defined as

$$\text{Frame}(\Pi \wedge \Sigma, \{\mathbf{x}\}) =_{df} \Pi \wedge \text{UnreachHeap}(\Pi \wedge \Sigma, \{\mathbf{x}\})$$

where $\text{UnreachHeap}(\Pi \wedge \Sigma, \{\mathbf{x}\})$ is the formula consisting of all *-conjuncts from Σ which are not in $\text{ReachHeap}(\Pi \wedge \Sigma, \{\mathbf{x}\})$.

Note that $\mathbf{x}=\mathbf{a} \wedge \mathbf{y}=\mathbf{b} \wedge \mathbf{z}=\mathbf{c}$ are used to record the snapshot of variables associated with U using the special variables \mathbf{a}, \mathbf{b} and \mathbf{c} . The symbolic execution of C_2 at line 8

```

Verify : AllSpec × P(SH × SH) × V → AllSpec ∪ {fail}
Algorithm Verify( $\mathcal{T}$ , SpecV, {C1; U(x;y); C2}(x0;y0))
1 SpecU := ∅
2 foreach (Pre, Post) ∈ SpecV do
3   S0 := ⟦C1⟧T{Pre ∧ y0=old(y0)}
4   if ⊤ ∈ S0 then return fail endif
5   foreach Δ1 ∈ S0 do
6     PreU := Local(Δ1, {x, y})
7     Denote z = fv(PreU) \ {x, y}
8     S := ⟦C2⟧TA{([old(b)/y] Frame(Δ1, {x, y}) ∧
        x=a ∧ y=b ∧ z=c, emp ∧ x=a ∧ y=b ∧ z=c)}
9     S' := { (Δ, M) | (Δ, M) ∈ S ∧ Δ ⊢ Post * true } ∪
        { (Δ * M', M * M') | (Δ, M) ∈ S ∧
        Δ ≠ Post * true ∧ Δ * [M'] ⊢ Post }
10    if ∃(Δ, M) ∈ S' . fv(M) ⊈ ReachVar(Δ, {a, b})
        then return (fail, M) endif
11    foreach (Δ, M) ∈ S' do
12      PreU := [a/x, b/y, c/z]PreU
13      PostU := sub_alias(M, {a, b, c})
14      g := (fv(PreU) ∩ fv(PostU)) ∪ {a, b}
15      SpecU := SpecU ∪ { (∃(fv(PreU) \ g) · PreU, PostU) }
16    end foreach
17  end foreach
18 end foreach
19 TU := CaseAnalysis( $\mathcal{T}$ , SpecU, U)
20 Post_Check(T ⊕ TU, SpecV, {C1; U; C2})
21 return T ⊕ TU
end Algorithm

```

Fig. 3. The main verification algorithm.

returns a set S of pairs (Δ, M) where Δ is a possible post state of C_2 and M records the discovered effect of U . At line 9, we check whether or not each Δ can establish the postcondition Post for the whole procedure. If not, another abduction $\Delta * [M'] \triangleright \text{Post}$ is invoked to discover further effect M' which is then added into M .

There can be some complication here. Note that there can be a potential bug in the program, or the given specification is not sufficient. As a consequence of that, the result M returned by our abductor may contain more information than what can be expected from U , in which case we cannot simply regard the whole M as the postcondition of U . For example, consider the code fragment `unknown(x); z:=y.next` with the precondition $x \mapsto \text{null}$. Before the second instruction (dereference of `y.next`) we use abduction to get $y \neq \text{null}$. However, noting the fact $y \notin \text{ReachVar}(\Delta, \{x\})$ where $\Delta = \text{emp} \wedge y \neq \text{null}$ is the state immediately after the `unknown` call plus the abduction result, we know that from the `unknown` call's parameters (x) , y is not reachable, and hence the `unknown` call will never establish a state where $y \neq \text{null}$. In that case we are assured that the procedure being verified cannot meet the specification.

To detect such a situation, we introduce the check in line 10. It tests whether the whole abduction result is reachable from variables accessed by U . If not, then the unreachable part cannot be expected from U , which indicates a possible bug in the program or some inconsistency between the program and its specification. In such cases, the algorithm returns an additional formula that can be used by a further analysis to either identify the bug or strengthen the specification. Recall the example presented in the previous paragraph: since $y \neq \text{null}$ cannot be established by the unknown call, if we add it to the precondition of the code fragment (to form a new precondition $x \mapsto \text{null} \wedge y \neq \text{null}$), then the verification with the new specification can move on and will potentially succeed. We will exemplify this later with experimental results.

The third step (lines 11–16) is to form the derived specifications for U in terms of variables \mathbf{a}, \mathbf{b} and g , where g denotes logical variables not directly accessed by U , but occurring in both pre- and postconditions. The formula $\text{sub_alias}(\mathbf{M}, \{\mathbf{a}, \mathbf{b}, \mathbf{c}\})$ is obtained from \mathbf{M} by replacing all variables with their aliases in $\{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$. It is defined as

$$\begin{aligned} \text{sub_alias}(\mathbf{M}, \{\mathbf{x}\}) =_{df} & (\{[x/x'] \mid x \in \{\mathbf{x}\} \wedge x' \in \text{aliases}(\mathbf{M}, x)\} \mathbf{M}) \wedge \\ & \wedge \{x = x' \mid x, x' \in \{\mathbf{x}\} \wedge x' \in \text{aliases}(\mathbf{M}, x)\} \end{aligned}$$

where a set of substitutions before a formula \mathbf{M} denotes the result of applying all those substitutions to \mathbf{M} , and $\text{aliases}(\mathbf{M}, x)$ returns all the aliases of x in \mathbf{M} .

Finally, at line 19, the obtained specifications Spec_U for U are passed to the case analysis algorithm (given in Figure 4) to derive the specifications of unknown procedures invoked in U . At line 20, we conduct a post analysis for soundness purpose.

2. Case analysis algorithm. In order to discover specifications for unknown procedures invoked in U , the algorithm in Figure 4 conducts a case analysis according to the structure of U . In the first case (line 2), U is simply an unknown call. In this situation, the algorithm simply returns all the pre-/postcondition pairs from Spec_U as the unknown procedure's specifications.

In the second case (line 4), U is an if construct and each branch contains an unknown block. The algorithm uses the main algorithm to verify the two branches separately with preconditions $\text{Pre} \wedge \mathbf{b}$ and $\text{Pre} \wedge \neg \mathbf{b}$ respectively, where Pre is one of the preconditions of the whole if. The results obtained from the two branches are then combined using the \uplus operator:

$$R_1 \uplus R_2 =_{df} \{(f, \text{Refine}(\text{Spec}_f^1 \cup \text{Spec}_f^2)) \mid (f, \text{Spec}_f^1) \in R_1 \wedge (f, \text{Spec}_f^2) \in R_2\}$$

where Refine is defined as

$$\begin{aligned} \text{Refine } (\emptyset) & =_{df} \emptyset \\ \text{Refine } (\{(Pre, Post)\} \cup \text{Spec}) & =_{df} \text{if } \exists (Pre', Post') \in \text{Spec} \cdot Pre' \preceq Pre \wedge Post' \preceq Post' \\ & \quad \text{then Refine}(\text{Spec}) \\ & \quad \text{else } \{(Pre, Post)\} \cup \text{Refine}(\text{Spec}) \end{aligned}$$

The intuition of Refine is to eliminate any specification $(Pre', Post')$ from a set if there exists a “stronger” one $(Pre, Post)$ such that $Pre \preceq Pre'$ and $Post' \preceq Post$. \uplus is to refine the union of two specification sets.

```

CaseAnalysis : AllSpec  $\times$   $\mathcal{P}(\text{SH} \times \text{SH}) \times U \rightarrow \text{AllSpec} \cup \{\text{fail}\}$ 
Algorithm CaseAnalysis( $\mathcal{T}, \text{Spec}_U, U$ )
1  switch  $U$ 
2  case unkFn( $\mathbf{x}; \mathbf{y}$ )
3    return  $\{(\text{unkFn}, \text{Spec}_U)\}$ 
4  case if  $b$  then  $V_1$  else  $V_2$  fi
5     $\text{Spec}_T := \{(\text{Pre} \wedge b, \text{Post}) \mid (\text{Pre}, \text{Post}) \in \text{Spec}_U\}$ 
6     $\text{Spec}_F := \{(\text{Pre} \wedge \neg b, \text{Post}) \mid (\text{Pre}, \text{Post}) \in \text{Spec}_U\}$ 
7     $R_1 := \text{Verify}(\mathcal{T}, \text{Spec}_T, V_1)$ 
8     $R_2 := \text{Verify}(\mathcal{T}, \text{Spec}_F, V_2)$ 
9    return  $R_1 \uplus R_2$ 
10 case if  $b$  then  $V$  else  $C$  fi
11    $\text{Spec}_T := \{(\text{Pre} \wedge b, \text{Post}) \mid (\text{Pre}, \text{Post}) \in \text{Spec}_U\}$ 
12    $R := \text{Verify}(\mathcal{T}, \text{Spec}_T, V)$ 
13   if  $\exists (\text{Pre}, \text{Post}) \in \text{Spec}_U, \Delta \in \llbracket C \rrbracket_{\mathcal{T}} \{ \text{Pre} \wedge \neg b \} \cdot$ 
         $\Delta = \top \vee \Delta \not\vdash \text{Post} * \text{true}$  then return fail
14   else return  $R$  endif
15 case if  $b$  then  $C$  else  $V$  fi (Similar to the previous case)
16 case while  $b$  do  $V$  od
17   Denote  $V$  as  $\{C_1; U_{(\mathbf{x}_1; \mathbf{y}_1)}; C_2\}_{(\mathbf{x}; \mathbf{y})}$ 
18   Define  $\text{loop}(\mathbf{x}; \mathbf{y}) \{ \text{if } b \text{ then } V; \text{loop}(\mathbf{x}; \mathbf{y}) \text{ fi} \}$ 
19    $\mathcal{T}' := \mathcal{T} \uplus \{ (\text{loop}, \text{Spec}_U) \}$ 
20   return  $\text{Verify}(\mathcal{T}', \text{Spec}_U, \{ \text{if } b \text{ then } V; \text{loop}(\mathbf{x}; \mathbf{y}) \text{ fi} \}_{(\mathbf{x}; \mathbf{y})})$ 
21 case unkFn( $\mathbf{x}_0; \mathbf{y}_0$ );  $\{ ; C_i; \text{unkFn}(\mathbf{x}_i; \mathbf{y}_i) \}_{i=1}^n$ 
22   return  $\{(\text{unkFn}, \text{SeqUnkCalls}(\mathcal{T}, \text{Spec}_U, U))\}$ 
end Algorithm

```

Fig. 4. The case analysis algorithm.

The third and fourth cases (line 10 and 15) are for if constructs which contain one unknown block in one branch. This is handled in a similar way as in the second case. The only difference is that, for the branch without unknown blocks, we need to verify it with the underlying semantics (line 13).

The fifth case is the while loop. In the motivating example in Section 2, we have shown that our approach is able to handle the verification of a tail-recursive function provided with both pre- and postconditions, our solution here is to translate the while loop into a tail-recursive function to verify it. As can be seen in lines 16 – 20, the algorithm generates a new function `loop` for the while loop, and takes the variables accessed by V to be its parameters. Note that the read-only variables (\mathbf{x}) become call-by-value parameters and other possibly mutable ones (\mathbf{y}) become call-by-reference parameters (which is one reason for us to introduce call-by-reference parameters in our language). The algorithm then adds the specifications found for the while loop as `loop`'s specifications into table \mathcal{T} , verifying it, and at the same time obtaining specifications for the unknown procedure in the while loop, using the main algorithm `Verify`.

In the last case (line 21), where U consists of multiple unknown procedure calls in sequence, the algorithm invokes another algorithm, `SeqUnkCalls`, to deal with it.

3. Verifying sequential unknown calls. To handle the most complicated case, unknown procedure calls in sequence, we still need the `SeqUnkCalls` algorithm. First we illustrate the brief idea using two sequential unknown procedure calls as an example, followed by the general algorithm.

Suppose we have

$$\{\text{Pre}\} \{\text{unkFn}_1(\mathbf{x}_0; \mathbf{y}_0); C; \text{unkFn}_2(\mathbf{x}_1; \mathbf{y}_1)\} \{\text{Post}\}$$

where C is the only known code fragment within the block. Our current solution attempts to find a common specification to capture both unknown procedures' behaviors.

The algorithm works in three steps. In the first step, it extracts the precondition for the first procedure, say Pre_U , from the given precondition Pre by extracting the part of heap that may be accessed by the call via \mathbf{x}_0 and \mathbf{y}_0 , which is similar to the first step of the main algorithm `Verify`. Aiming at a general specification for both unknown calls, it then assumes that the second procedure has a similar precondition Pre_U . In the second step, it symbolically executes the code fragment C with the help of the abductor, to discover a crude postcondition, say Post_U^0 , expected from the first unknown call. This is similar to the second step of the main algorithm `Verify`, except that the postcondition for C is now assumed to be Pre_U . In the third step, the algorithm takes Post_U^0 (with appropriate variable substitutions) as the postcondition of the second unknown call, and checks whether or not the derived post (Post_U^0) satisfies Post . If not, it invokes another abduction to strengthen Post_U^0 to obtain the final postcondition Post_U for the unknown procedures. Note that this strengthening does not affect soundness: the strengthened Post_U can still be used as a general postcondition for both unknown procedures.

Figure 5 presents the algorithm to infer specifications for n ($n \geq 2$) unknown calls in sequence. As aforementioned, given a block of $(n+1)$ unknown procedure calls with n pieces of known code blocks sandwiched among them $(\text{unkFn}(\mathbf{x}_0; \mathbf{y}_0) \{; C_i; \text{unkFn}(\mathbf{x}_i; \mathbf{y}_i)\}_{i=1}^n$ in line 1), and the specification $(\text{Pre}, \text{Post})$ (line 3) for such a block, our approach generally works in three steps: first, to compute a precondition for the unknown calls; second, to verify each code fragment C_i ($i = 1, \dots, n$) with abduction to collect expected behavior of the unknown calls (as part of their postcondition); third, to guarantee that the collected postcondition satisfies Post . If not, then another abduction is conducted to strengthen the gained postcondition to ensure this.

The first step is completed by lines 4 to 6. The local part of Pre is extracted w.r.t. the first unknown call's parameters \mathbf{x}_0 and \mathbf{y}_0 . Other free variables are distinguished as \mathbf{z}_0 , which may be ghost variables. Finally the precondition is found in terms of special logical variables \mathbf{a} , \mathbf{b} and \mathbf{c} .

The second step is performed over each $C_i; \text{unkFn}(\mathbf{x}_i; \mathbf{y}_i)$. Its main idea is to take the postcondition generated for the last unknown call (Post_{i-1}), plus the frame part during the entailment check against Pre_{i-1} , as the post-state of $\text{unkFn}(\mathbf{x}_{i-1}; \mathbf{y}_{i-1})$, and try to verify C_i beginning with such a state, using abduction when necessary (line 9). After the verification we get S_i containing abstract states before $\text{unkFn}(\mathbf{x}_i; \mathbf{y}_i)$, and we want those states to satisfy its precondition Pre_U subject to substitution. Note that during the verification of C_i and the last satisfaction checking we may use abduction to strengthen the program state, whose results reflect the expected behavior of $\text{unkFn}(\mathbf{x}_{i-1}; \mathbf{y}_{i-1})$ and are accumulated as its expected postcondition. Hence we achieve a sufficiently strong postcondition for each unknown call.

```

SeqUnkCalls : AllSpec ×  $\mathcal{P}(\text{SH} \times \text{SH}) \times U \rightarrow \mathcal{P}(\text{SH} \times (\text{SH} \cup \{\top\}))$ 
Algorithm SeqUnkCalls( $\mathcal{T}, \text{Spec}_U, U$ )
1  Denote  $U$  as  $\text{unkFn}(\mathbf{x}_0; \mathbf{y}_0) \{; C_i; \text{unkFn}(\mathbf{x}_i; \mathbf{y}_i)\}_{i=1}^n$ 
2   $R := \emptyset$ 
3  foreach (Pre, Post)  $\in \text{Spec}_U$  do
4     $\text{Pre}_U := \text{Local}(\text{Pre}, \{\mathbf{x}_0, \mathbf{y}_0\})$ 
5    Denote  $\mathbf{z}_0 = \text{fv}(\text{Pre}_U) \setminus \{\mathbf{x}_0, \mathbf{y}_0\}$ 
6     $\text{Pre}_U := [\mathbf{a}/\mathbf{x}_0, \mathbf{b}/\mathbf{y}_0, \mathbf{c}/\mathbf{z}_0]\text{Pre}_U$ 
7     $S'_0 := \{(\text{Pre} \wedge \mathbf{y}_0 = \text{old}(\mathbf{y}_0), \text{emp} \wedge \mathbf{a} = \mathbf{x}_0 \wedge \mathbf{b} = \mathbf{y}_0 \wedge \mathbf{c} = \mathbf{z}_0)\}$ 
8    for  $i := 1$  to  $n$  do
9       $S_i := \llbracket C_i \rrbracket_{\mathcal{T}}^{\wedge} \{ (\text{Post}_{i-1} * [\text{old}(\mathbf{b})/\mathbf{y}_{i-1}] \text{Frame}(\Delta_{i-1}, \{\mathbf{x}_{i-1}, \mathbf{y}_{i-1}\}), \text{Post}_{i-1}) \mid$ 
       $(\Delta_{i-1}, M_{i-1}) \in S'_{i-1} \wedge \text{Post}_{i-1} = ([\mathbf{x}_{i-1}/\mathbf{a}, \mathbf{y}_{i-1}/\mathbf{b}, \mathbf{z}_{i-1}/\mathbf{c}] \text{sub\_alias}(\mathbf{M}_{i-1}, \{\mathbf{a}, \mathbf{b}, \mathbf{c}\})) \wedge \mathbf{a} = \mathbf{x}_{i-1} \wedge \mathbf{b} = \mathbf{y}_{i-1} \wedge \mathbf{c} = \mathbf{z}_{i-1} \}$  where  $\mathbf{z}_{i-1}$  is fresh
10      $S'_i := \{(\Delta, M) \mid (\Delta, M) \in S_i \wedge \sigma \Delta \vdash \exists \mathbf{c} \cdot \text{Pre}_U * \text{true}\} \cup \{(\Delta * M', M * M') \mid$ 
       $(\Delta, M) \in S_i \wedge \sigma \Delta \not\vdash \exists \mathbf{c} \cdot \text{Pre}_U * \text{true} \wedge \sigma \Delta * [M'] \triangleright \exists \mathbf{c} \cdot \text{Pre}_U\}$ 
      where  $\sigma = [\mathbf{a}/\mathbf{x}_i, \mathbf{b}/\mathbf{y}_i]$ 
11     if  $\exists (\Delta, M) \in S'_i \cdot \text{fv}(M) \not\subseteq \text{ReachVar}(\Delta, \{\mathbf{a}, \mathbf{b}\})$  then
      return (fail, Local( $M, \{\mathbf{x}_0, \text{old}(\mathbf{y}_0)\}$ )) endif
12    end for
13     $S_{n+1} := \{ (\text{Post}_n * [\text{old}(\mathbf{b})/\mathbf{y}_n] \text{Frame}(\Delta_n, \{\mathbf{x}_n, \mathbf{y}_n\}), \text{Post}_n) \mid (\Delta_n, M_n) \in S'_n \wedge$ 
       $\text{Post}_n = ([\mathbf{x}_n/\mathbf{a}, \mathbf{y}_n/\mathbf{b}, \mathbf{z}_n/\mathbf{c}] \text{sub\_alias}(M_n, \{\mathbf{a}, \mathbf{b}, \mathbf{c}\})) \wedge$ 
       $\mathbf{a} = \mathbf{x}_n \wedge \mathbf{b} = \mathbf{y}_n \wedge \mathbf{c} = \mathbf{z}_n \}$  where  $\mathbf{z}_n$  is fresh
14     $S'_{n+1} := \{(\Delta, M) \mid (\Delta, M) \in S_{n+1} \wedge \Delta \vdash \text{Post} * \text{true}\} \cup \{(\Delta * M', M * M') \mid$ 
       $(\Delta, M) \in S_{n+1} \wedge \Delta \not\vdash \text{Post} * \text{true} \wedge \Delta * [M'] \triangleright \text{Post}\}$ 
15    if  $\exists (\Delta, M) \in S'_{n+1} \cdot \text{fv}(M) \not\subseteq \text{ReachVar}(\Delta, \{\mathbf{a}, \mathbf{b}\})$  then
      return (fail, Local( $M, \{\mathbf{x}_0, \text{old}(\mathbf{y}_0)\}$ )) endif
16    foreach  $(\Delta, M) \in S'_{n+1}$  do
17       $\text{Post}_U := \text{sub\_alias}(M, \{\mathbf{a}, \mathbf{b}, \mathbf{c}\})$ 
18       $g := \text{fv}(\text{Pre}_U) \cap \text{fv}(\text{Post}_U) \setminus \{\mathbf{a}, \mathbf{b}\}$ 
19       $R := \text{Refine}(R \cup \{(\exists \text{fv}(\text{Pre}_U) \setminus (g \cup \{\mathbf{a}, \mathbf{b}\}) \cdot \text{Pre}_U, \text{Post}_U)\})$ 
20    end foreach
21  end foreach
22  return  $R$ 
end Algorithm

```

Fig. 5. Algorithm for sequential unknown calls.

The third step is similar to the first algorithm: it checks whether the final abstract state entails the postcondition of the whole block, and strengthens the final abstract state with abduction if it cannot. Then the ghost variables are recognized and processed analogously to the first algorithm. Finally the strongest specifications discovered for those unknown procedures are returned.

Note that our current solution tries to find a common specification (Pre, Post) suitable for all the unknown procedures. Generally we may allow the unknown procedures to have different specifications. In theory, this can be achieved by a more in-depth analysis

which examines the known code fragments in between those unknown calls. That is, by analyzing the code fragment C_i we would hopefully obtain a postcondition for the $(i-1)$ -th procedure and a precondition for the i -th. In the case of two unknown calls $\text{unkFn}_0(\mathbf{x}_0; \mathbf{y}_0); C_1; \text{unkFn}_1(\mathbf{x}_1; \mathbf{y}_1)$, the precondition for unkFn_0 and the postcondition for unkFn_1 can be obtained as usual (by analyzing the code before unkFn_0 and after unkFn_1 respectively). To derive the postcondition Post_0 for unkFn_0 and the precondition Pre_1 for unkFn_1 , we initialize Post_0 to be emp to start a forward analysis over C_1 with abduction, to accumulate (via abduction) the expected behavior of unkFn_0 (for C_1 to be verified) as Post_0 , and extract a formula (which is relevant to the footprint of unkFn_1) from the abstract state at the end of C_1 as Pre_1 . However, our initial experiments show that, unless the fragment C_1 is sufficiently complex to expose enough information expected from unkFn_0 , the derived Post_0 and Pre_1 can be rather weak. As a consequence, the derived specification for unkFn_0 can be too weak (with a weak postcondition) and the one for unkFn_1 can be too strong (with a weak precondition). It remains an open problem how we might tune the derived results to obtain more reasonable specifications. We conjecture that certain heuristics might help and we will explore this further in our future work.

6.2. Soundness and Termination

Informally, in the presence of unknown procedure calls, the soundness of the verification signifies that, a program is successfully verified against its specifications, if all the unknown procedures that it invokes conform to the specifications discovered by the verification algorithm. Therefore, the correctness of the program depends on a (possible) further verification for the unknown procedures.

The soundness of our verification algorithm is guaranteed at line 20 in the algorithm `Verify`. Upon return, the algorithm conducts another forward verification on the whole program, assuming that the unknown procedures are already verified against the discovered specifications. As in Calcagno et al. [6], this final check rules out any potentially unsound pre/post pairs (due to the use of abstraction); therefore, it ensures the soundness of our verification.

The termination of our verification algorithms is based on two facts: the finiteness of verification input (program and its specifications), and the termination of our semantics. Firstly, all the loops and (recursive) algorithm invocations in the three verification algorithms are performed over the input program and specifications to be verified, where each time of iteration deals with part of the specifications/states and each invocation processes part of the input program structurally. Therefore, as long as all such processing terminates, the whole verification will terminate. This fact is further guaranteed with the termination of our abstract semantics, which can be proved by claiming the finiteness of program and logical variables, and hence the finiteness of all possible abstract states [9]. Hence we can conclude that our verification algorithm terminates with results of either fail or successfully returned specifications that the unknown procedures must conform.

7. Experiments and Evaluation

We have implemented the two abstract semantics and the verification algorithms with Objective Caml and evaluated them over some list processing programs to test their

viability and precision. The results are in Table 1. The first and second columns denote the programs used for evaluation [5, 8] and their time consumption, respectively. We manually hide some code in the original programs as calls to unknown procedures, for which we try to discover specifications during the verification process. The third column shows the partial specification for each program that is provided as input to our algorithm. The last column exhibits the discovered specifications for unknown procedures inside those programs. The programs are listed in Figure 6.

Program	Time(s)	Program Specification	Discovered Unknown Specification
findLast (<i>x</i> ; <i>z</i>)	0.00135	Pre := $\text{ls}(x, \text{null}) \wedge x \neq \text{null}$ Post := $\text{ls}(x, z) * z \mapsto \text{null}$	Pre := $\exists w \cdot a \mapsto w * \text{ls}(w, \text{null}) \wedge$ $a \neq \text{null} \wedge w \neq \text{null}$ Post := $\text{ls}(a, b) * \text{ls}(b, \text{null}) \wedge b \neq \text{null}$
append (<i>y</i> ; <i>x</i>)	0.00216	Pre := $\text{ls}(x, \text{null}) * \text{ls}(y, \text{null})$ Post := $\text{ls}(x, y) * \text{ls}(y, \text{null})$	Pre := $\text{ls}(a, \text{null})$ Post := $\text{ls}(a, b) * \text{ls}(b, \text{null})$
copy (<i>x</i> ; <i>y</i>)	0.00204	Pre := $\text{ls}(x, \text{null})$ Post := $\text{ls}(x, \text{null}) * \text{ls}(y, \text{null})$	Pre := $\text{ls}(a, \text{null})$ Post := $\text{ls}(a, \text{null})$
revCopy (<i>x</i> ; <i>y</i>)	0.00107	Pre := $\text{ls}(x, \text{null})$ Post := $\text{ls}(x, \text{null}) * \text{ls}(y, \text{old}(y))$	Pre := true Post := true
clear (; <i>x</i>)	0.00239	Pre := $\text{ls}(x, \text{null})$ Post := $\text{emp} \wedge x = \text{null}$	Pre := $a \mapsto b * \text{ls}(b, \text{null})$ Post := $\text{ls}(b, \text{null})$
appendThree (<i>y</i> , <i>z</i> ; <i>x</i>)	0.00315	Pre := $\text{ls}(x, \text{null}) * \text{ls}(y, \text{null}) * \text{ls}(z, \text{null})$ Post := $\text{ls}(x, \text{null})$	Pre := $\text{ls}(a, \text{null}) * \text{ls}(b, \text{null})$ Post := $\text{ls}(a, \text{null})$
towardsLast (<i>x</i> ; <i>y</i>)	0.00428	Pre := $\text{ls}(x, \text{null}) \wedge x \neq \text{null}$ Post := $\text{ls}(x, y) * \text{ls}(y, \text{null}) \wedge x \neq \text{null}$	Pre := $\text{ls}(a, \text{null}) \wedge a \neq \text{null}$ Post := $\text{ls}(a, \text{null}) * \text{ls}(b, \text{null}) \wedge a \neq \text{null} \wedge b \neq \text{null}$
iWillFail (<i>x</i> , <i>y</i> ;)	0.00066	Pre := $\text{ls}(x, \text{null})$ Post := $\text{ls}(x, \text{null}) * \text{ls}(y, \text{null})$	(fail, $\text{ls}(y, \text{null})$)

Table 1. Experimental results for simple list processing programs.

Here we note down two observations on the experimental results. The first is that the discovered specifications for the unknown procedures are more general than we would have expected. Bear in mind that we have replaced some code from those programs with unknown calls. We have compared the inferred specifications for those unknown calls with the original code. The results show that the specifications derived by our algorithm not only fully capture the behaviors of the replaced code, but also suggest other possible implementations. A case in point is our motivating example **findLast** given in Sec 2, where the original code replaced by unknown call is $y := x.\text{next}$ with the post-state $x \mapsto y * \text{ls}(y, \text{null}) \wedge y \neq \text{null}$. According to our result, the post-state can be more general, namely, $\text{ls}(x, y) * \text{ls}(y, \text{null}) \wedge y \neq \text{null}$. This suggests that as long as y is not **null**, the unknown call can traverse any number of nodes towards the tail of the list.

The second observation is about the last program listed in the table, a procedure called **iWillFail**, whose code is given as follows:

```
void iWillFail(node x, node y; ) { unknown(x; ) }
```

```

findLast(x; ref z) {
  node w := [x], y;
  if w = null then z := x
  else
    unknown(x; y);
    findLast(y; z)
  fi
}

copy(x; ref y) {
  if x = null then y := null
  else
    w := [x];
    copy(w; y);
    unknown(; y);
  fi
}

clear(; ref x) {
  if x = null then skip
  else
    w := [x];
    unknown(x, w; );
    x := w;
    clear(; x)
  fi
}

towardsLast(x; y) {
  unknown(; x);
  unknown(x; y)
}

append(y; ref x) {
  if x = null then x := y
  else
    unknown(x; w);
    if w = null then [x] := y
    else append(y; w)
  fi
}

revCopy(x; ref y) {
  if x = null then skip
  else
    unknown(; y);
    w := [x];
    revCopy(w; y)
  fi
}

appendThree(y, z; ref x) {
  unknown(x, y; );
  unknown(x, z; )
}

iWillFail(x, y; ) {
  unknown(x; )
}

```

Fig. 6. Code of experimental examples.

This program is expected to be verified against the specification ($\text{Pre} = \text{ls}(x, \text{null})$, $\text{Post} = \text{ls}(x, \text{null}) * \text{ls}(y, \text{null})$). Our verification fails and returns an additional formula $\text{ls}(y, \text{null})$ from our abduction process. A further analysis reveals that the failure is actually due to the given specification where the precondition Pre is too weak for the program to establish the postcondition Post : since y is not reachable from the parameters of the `unknown` call, no implementation of the `unknown` call can establish the postcondition Post involving y . The returned formula from our verification can then be used to strengthen the given specification. In this case, if we add $\text{ls}(y, \text{null})$ into Pre via separation conjunction, the verification will succeed with the specification ($\text{Pre}_u = \text{ls}(x, \text{null})$, $\text{Post}_u = \text{ls}(x, \text{null})$) discovered for the procedure `unknown`.

We have also conducted some experiments to test our approach’s performance and scalability. The results are shown in Table 2.

Program	Lines of code	Time(s)	Memory(Mb)	Specs discovered
<code>list.h list.c</code>	474	0.2	14.54	4
<code>tasks.h tasks.c</code>	3150	11.2	17.972	4
<code>wd-stat.h wd-stat.c</code> <code>ctrins.h ctrins.c</code>	1404	1.232	16.664	10
<code>kr-db.h kr-db.c</code>	1674	1.837	16.636	12
<code>krbid.h krbid.c</code> <code>krpage.h krpage.c</code>	2873	6.44	16.676	11
<code>kdbadapt.h kdbadapt.c</code> <code>kdbview.h kdbview.c</code>	3208	7.88	17.54	13

Table 2. Experimental results for performance.

These results were achieved with an Intel Core 2 Quad CPU 2.66GHz with 8Gb memory. The verified programs are either source code from FreeRTOS [2] (`list.c` and `tasks.c`), or some working code created by the author (`wd-stat.c`, `ctrins.c`, `kr-db.c`, `krbid.c`, `krpage.c`, `kdbadapt.c` and `kdbview.c`). All these programs mainly deal with pointer-based linear data structures, such as singly-linked and/or doubly-linked lists. For example, `list.c` provides functions to initialize and modify lists, while `tasks.c` calls those functions to manipulate several lists during runtime. The other programs, `wd-stat.c`, `ctrins.c`, `kr-db.c`, `krbid.c`, `krpage.c`, `kdbadapt.c` and `kdbview.c` maintain some list-based vectors for runtime data storage, or hashtables implemented with a series of linked lists. Meanwhile, these programs have many function calls in them, like invocation of library functions or calls of other functions in the project.

As for verification purpose, we provide each function to be verified with its specification and run our algorithm over it. The unknown functions are manually assigned, such as library function calls which mainly consist of memory allocation functions like `malloc`, and functions that reside in other programs of the project. Since we only consider pointer safety of linked data structures at the moment, we did some modification to the programs to revise the code not relevant to our verification (like data types `int` and `char` for instance), such that the programs can be successfully verified.

The experimental results suggest that our approach might be able to scale up as a verification system for pointer safety, albeit we are trying to do more experiments, as well as optimize the algorithm, to justify this. What is more, we have tested the top-down feature of our abductive based verification system in the experiments. For instance, the specifications gained for “unknown” functions being invoked in `tasks.c` cohere with the ones for the corresponding callee functions in `list.c`. The same situation also applies to other test cases. Based on such results, we will investigate more towards this direction to uncover the power of a top-down abductive-based system as a viable alternative to the current bottom-up approaches [6, 13].

8. Conclusion

It is a challenging problem to automatically verify (even pointer safety of) heap-manipulating imperative programs with unknown procedure calls. We propose a novel approach to this problem, which infers expected specifications for unknown procedures from their calling contexts during the verification process. The program is proven correct subject to the condition that the invoked unknown procedures meet the inferred specifications. We employ a forward shape analysis with separation logic and an enhanced abductive reasoning mechanism to synthesize both pre- and postconditions of the unknown procedure. As a proof of concept, we have also implemented a prototype system to test the viability of the proposed approach.

There are two possible future directions. One is to explore a more general solution for unknown calls in sequence as discussed, e.g., it might be possible for us to invent some heuristics to strengthen the postcondition for the first unknown call, so that the precondition of the second can be strengthened accordingly, to achieve more reasonable specifications for both. Another direction is to extend this method to an abstract domain combining separation and numerical information [7], so that more general properties, such as memory safety and functional correctness, can be specified and verified. We envisage that, with the combined domain, the abstract semantics and analysis algorithms will remain conceivably the same, but the abduction will be redefined to discover the anti-frames for the newly introduced numerical features.

Acknowledgement

This work was supported in part by the EPSRC (Engineering and Physical Sciences Research Council) projects [EP/E021948/1 and EP/G042322/1]. We thank Hongseok Yang for his encouragement and invaluable comments and Dino Distefano for very useful discussions. We also acknowledge the precious insights from the anonymous reviewers.

References

- [1] Glenn Ammons, Rastislav Bodik, and James R. Larus. Mining specifications. In *Proc. of 29th Principles of Programming Languages (POPL)*, New York, NY, USA, pages 4–16. ACM Press, 2002.
- [2] Richard Barry. *FreeRTOS Reference Manual — API Functions and Configuration Options*. Real Time Engineers Ltd, 2009.
- [3] Josh Berdine, Cristiano Calcagno, and Peter W. O’Hearn. A decidable fragment of separation logic. In *Proc. of Foundations of Software Technology and Theoretical Computer Science (FSTTCS)*, Chennai, India, LNCS 3328, pages 97–109. Springer, December 2004.
- [4] Josh Berdine, Cristiano Calcagno, and Peter W. O’Hearn. Symbolic execution with separation logic. In *Proc. of 3rd Asian Symposium on Programming Languages and Systems (APLAS)*, Tsukuba, Japan, LNCS 3780, pages 52–68. Springer, 2005.
- [5] Cristiano Calcagno, Dino Distefano, Peter W. O’Hearn, and Hongseok Yang. Footprint analysis: A shape analysis that discovers preconditions. In *Proc. of 14th International Symposium on Static Analysis (SAS)*, Denmark, LNCS 4634, pages 402–418. Springer, 2007.
- [6] Cristiano Calcagno, Dino Distefano, Peter W. O’Hearn, and Hongseok Yang. Compositional shape analysis by means of bi-abduction. In *Proc. of 36th Principles of Programming Languages (POPL)*, Savannah, Georgia, USA, pages 289–300. ACM Press, January 2009.

- [7] Wei-Ngan Chin, Cristina David, Huu Hai Nguyen, and Shengchao Qin. Automated verification of shape, size and bag properties. In *Proc. of 12th International Conference on Engineering of Complex Computer Systems (ICECCS)*, Auckland, New Zealand, pages 307–320. IEEE, 2007.
- [8] Florin Craciun, Chenguang Luo, Guanhua He, Shengchao Qin, and Wei-Ngan Chin. Discovering specifications for unknown procedures (work in progress). In *Proc. of Workshop of Invariant Generation (WING)*, pages 30–44, March 2009.
- [9] Dino Distefano, Peter W. O’Hearn, and Hongseok Yang. A local shape analysis based on separation logic. In *Proc. of 12th Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*, Vienna, Austria, LNCS 3920, pages 287–302. Springer, 2006.
- [10] Maryam Emami, Rakesh Ghiya, and Laurie J. Hendren. Context-sensitive interprocedural points-to analysis in the presence of function pointers. In *Proc. of Programming Language Design and Implementation (PLDI)*, New York, NY, USA, pages 242–256. ACM Press, 1994.
- [11] Roberto Giacobazzi. Abductive analysis of modular logic programs. In *Proc. of the International Symposium on Logic Programming (ISLP)*, pages 377–391. The MIT Press, 1994.
- [12] Denis Gopan and Thomas Reps. Low-level library analysis and summarization. In *Proc. of 19th Computer Aided Verification (CAV)*, Berlin, Germany, LNCS 4590, pages 68–81. Springer, 2007.
- [13] Bhargav S. Gulavani, Supratik Chakraborty, Ganesan Ramalingam, and Aditya V. Nori. Bottom-up shape analysis. In *Proc. of 16th International Symposium on Static Analysis (SAS)*, Los Angeles, CA, USA, LNCS 5673, pages 188–204. Springer, 2009.
- [14] Cliff Jones, Peter W. O’Hearn, and Jim Woodcock. Verified Software: A Grand Challenge. *IEEE Computer*, 39(4):93–95, April 2006.
- [15] Peter Müller, Arnd Poetzsch-Heffter, and Gary T. Leavens. Modular specification of frame properties in JML. In *Concurrency and Computation: Practice and Experience*, 15:117–154, 2003.
- [16] Huu Hai Nguyen, Cristina David, Shengchao Qin, and Wei-Ngan Chin. Automated verification of shape and size properties via separation logic. In *Proc. of 8th Verification, Model Checking, and Abstract Interpretation (VMCAI)*, Nice, France, LNCS 4349, pages 251–266. Springer, 2007.
- [17] Flemming Nielson, Hanne Riis Nielson, and Chris Hankin. *Principles of Program Analysis*. Springer, 2005.
- [18] Peter W. O’Hearn, Hongseok Yang, and John C. Reynolds. Separation and information hiding. In *Proc. of 31st Principles of Programming Languages (POPL)*, Venice, Italy, pages 268–280. ACM Press, January 2004.
- [19] John C. Reynolds. Separation logic: a logic for shared mutable data structures. In *Proc. of 17th IEEE Symposium on Logic in Computer Science (LICS)*, Copenhagen, Denmark, pages 55–74. IEEE, 2002.
- [20] Roger Sessions. *COM and DCOM: Microsoft’s vision for distributed objects*. John Wiley & Sons, Inc., New York, NY, USA, 1998.
- [21] Jim Woodcock. Verified software grand challenge. In *Proc. of 14th International Symposium on Formal Methods (FM)*, Hamilton, Canada, LNCS 4085, pages 617–617. Springer, 2006.
- [22] Hongseok Yang, Oukse Lee, Josh Berdine, Cristiano Calcagno, Byron Cook, Dino Distefano, and Peter W. O’Hearn. Scalable shape analysis for systems code. In *Proc. of 20th Computer Aided Verification (CAV)*, Princeton, NJ, USA, LNCS 5123, pages 385–398. Springer, April 2008.