

# Automatically Refining Partial Specifications for Program Verification

Shengchao Qin<sup>1</sup>, Chenguang Luo<sup>2\*</sup>, Wei-Ngan Chin<sup>3</sup>, and Guanhua He<sup>4</sup>

<sup>1</sup>Teesside University

<sup>2</sup>Citigroup Inc.

<sup>3</sup>National University of Singapore

<sup>4</sup>Durham University

**Abstract.** Automatically verifying heap-manipulating programs is a challenging task, especially when dealing with complex data structures with strong invariants, such as sorted lists and AVL/red-black trees. The verification process can greatly benefit from human assistance through specification annotations, but this process requires much intellectual effort from users and is error-prone. In this paper, we propose a new approach to program verification that allows users to provide only partial specification to methods. Our approach will then refine the given annotation into a more complete specification by discovering missing constraints. The discovered constraints may involve both numerical and multi-set properties that could be later confirmed or revised by users. We further augment our approach by requiring only partial specification to be given for primary methods. Specifications for loops and auxiliary methods can then be systematically discovered by our augmented mechanism, with the help of information propagated from the primary methods. Our work is aimed at verifying beyond shape properties, with the eventual goal of analysing full functional properties for pointer-based data structures. Initial experiments have confirmed that we can automatically refine partial specifications with non-trivial constraints, thus making it easier for users to handle specifications with richer properties.

## 1 Introduction

Human assistance is often essential in (semi-) automated program verification. The user may supply annotations at certain program point, such as loop invariants and/or method specifications. These annotations can greatly narrow down the possible program states at that point, and avoid fixed-point calculation which could be expensive and may be less precise than the user's insight.

However, an obvious disadvantage of user annotation concerns its scalability, since programs to be analysed may be complicated and their functions are also diverse. Therefore, it is not preferable to require the user to provide specification for each method and invariant for each loop when verifying a relatively large software system. Meanwhile, human is liable to make mistakes. A programmer may under-specify with too weak a precondition or over-specify with too strong a postcondition. Such mistakes could lead to failed verification, and it may be

---

\* Work done when the author was affiliated with Durham University.

difficult for the user to discover whether the error is due to a real bug in the program, or an inappropriately supplied annotation.

To balance verification quality and human effort, we provide a novel approach to the verification of heap manipulating programs, which has long since been a challenging problem. To deal with such programs, which manipulate heap-allocated shared mutable data structures, one needs to keep track of not only “shape” information (for deep heap properties) but also related “pure” properties, such as structural numerical information (size and height), relational numerical information (balanced and sortedness properties), and content information (multi-set of symbolic values). Under our framework, the user is expected to provide partial specifications for *primary* methods with only *shape* information. Our verification will then take over the rest of the work to refine those partial specifications with derived (pure) constraints which should be satisfied by the program, or report a possible program bug if the given specifications are rejected by our verifier. This is more beneficial compared with previous works [34, 35], where users must provide full specifications for each method and invariants for each loop. This is also significantly different from the compositional shape analysis [5, 15, 45]. In spite of a higher level of automation, their analysis focuses on pointer safety only and deals with a few built-in predicates over the shape domain. Our work targets at both memory safety and functional correctness and supports user-defined predicates over several abstract domains (such as shape, numerical, multi-set).

Our approach allows the user to design their predicates for shapes and relative properties, to capture the desired level of program correctness to be verified. For example, with a singly-linked list structure `data node { int val; node next; }`, a user interested in pointer-safety may define a list shape predicate (as in [5, 15]):

$$\text{list}(p) \equiv (p = \text{null}) \vee (\exists i, q. p \mapsto \text{node}(i, q) * \text{list}(q))$$

Note that in the inductive case, the separation conjunction  $*$  ([40]) ensures that two heap portions (the head node and the tail list) are domain-disjoint.

Yet another user may be interested to track also the length of a list to analyse quantitative measures, such as heap/stack resource usage, using

$$\text{ll}(p, n) \equiv (p = \text{null} \wedge n = 0) \vee (p \mapsto \text{node}(\_, q) * \text{ll}(q, m) \wedge n = m + 1)$$

Note that unbound variables, such as  $q$  and  $m$ , are implicitly existentially quantified, and  $\_$  is used to denote an existentially quantified anonymous variable. This predicate may be extended to capture the content information, to support a higher-level of correctness with multi-set (bag) property:

$$\text{llB}(p, S) \equiv (p = \text{null} \wedge S = \emptyset) \vee (p \mapsto \text{node}(v, q) * \text{llB}(q, S_1) \wedge S = \{v\} \sqcup S_1)$$

where the length of the list is implicitly captured by the cardinality  $|S|$ . A further strengthening can capture also the sortedness property:

$$\text{sllB}(p, S) \equiv (p = \text{null} \wedge S = \emptyset) \vee (p \mapsto \text{node}(v, q) * \text{sllB}(q, S_1) \wedge S = \{v\} \sqcup S_1 \wedge (\forall x \in S_1. v \leq x))$$

Therefore, the user can provide predicate definitions w.r.t. their required correctness level and program properties. These predicates may be non-trivial but can be reused multiple times for specifications of different methods. We have also built a library of predicates with respect to commonly-used data structures and useful program properties.

Based on these predicates, the user is expected to provide partial specifications for some primary methods which are the main objects of verification. Say, for a sorting algorithm taking  $x$  as input parameter that is expected to be non-null, the user may provide  $llB(x, S_1)$  as precondition and  $sllB(x, S_2)$  as postcondition, and our approach will refine the specification as  $llB(x, S_1) \wedge x \neq null$  for pre, and  $sllB(x, S_2) \wedge S_1 = S_2$  for post. Here we need user annotations as the initial specification, because we reserve the flexibility of verification w.r.t. different program properties at various correctness levels. For example, our approach can also verify the same algorithm, but for the following refined specifications:

```

requires list(x)  ∧ x≠null ensures list(x)
requires ll(x, n1) ∧ n1>0 ensures ll(x, n2) ∧ n1=n2
requires llB(x, S1) ∧ x≠null ensures llB(x, S2) ∧ S1=S2
requires llB(x, S) ∧ x≠null ensures ll(x, n) ∧ |S|=n

```

where the discovered missing constraints are shown in shaded form. The last pre/post can be omitted in our approach if we are given a coercion lemma [34] between  $x::ll\langle n \rangle$  and  $x::llB\langle S \rangle$ . This can help reduce the number of redundant specifications considered (or synthesised for auxiliary methods).

To summarise, our proposal for refining partial specification is aimed at harnessing the synergy between human's insights and machine's capability at automated program analysis. In particular, human's guidance can help narrow down on the most important of the numerous specifications that are possible with each program code, while automation by machine is important for minimising on the tedium faced by users. Our proposal has the following characteristics:

- *Specification completion*: This verification refines the specification from three aspects, namely, the constraints needed in the precondition for memory and code safety, the constraints in postcondition to link the method's pre- and post-states, and the constraints that the method's post-state satisfies.
- *Flexibility*: We allow the user to define their own predicates for the program properties they want to verify, so as to provide different levels of correctness. Meanwhile we aim at, and have covered much of, full functional correctness of pointer-manipulating programs such as data structure shapes, pointer safety, structural/relational numerical constraints, and bag information.
- *Reduction of user annotations*: Our approach uses program analysis techniques effectively to reduce users' annotations. As for our experiments, the user only has to supply the partial specifications for primary methods, and the analysis will compute pre- and postconditions for loops and auxiliary methods as well as refine primary methods' specifications.

- *Semi-Automation*: We classify our approach as semi-automatic, because the user is allowed to interfere and guide the verification at any point. For instance, they may provide invariant for a loop instead of our automated invariant generation, or choose some other constraints as refinement from what the verification has discovered.

We have built a prototype implementation and carried out a number of experiments to confirm the viability of the approach as described in Section 5. In what follows, we will first depict our approach informally using a motivating example and present technical details thereafter. More related works and concluding remarks come after the experimental results.

## 2 The Approach

In this section, we briefly introduce the HIP/SLEEK system as the base of our verification and refinement. We then use some motivating examples to informally illustrate our approach.

### 2.1 The Hip/Sleek System

Separation logic [24, 40] extends Hoare logic to support reasoning about shared mutable data structures. It adds two more connectives to classical logic: separation conjunction  $*$  and spacial implication  $-*$ . The formula  $p_1 * p_2$  asserts that two heaps described by  $p_1$  and  $p_2$  are domain-disjoint, while  $p_1 -* p_2$  asserts that if the current heap is extended with a disjoint heap described by  $p_1$ , then  $p_2$  holds in the extended heap. In this paper we only use separation conjunction.

For better flexibility and expressivity, HIP/SLEEK allows users to define inductive shape predicates to leverage both shape and pure properties. We have illustrated several of these shape predicate definitions in the last section. For more involved examples, based on a data structure definition `data node2 { int val; node2 prev; node2 next; }`, one may define the predicate below to specify sorted doubly-linked list segments:

$$\text{sdlB}\langle p, q, S \rangle \equiv (\text{root} = q \wedge S = \emptyset) \vee (\text{root}::\text{node2}\langle v, p, r \rangle * r::\text{sdlB}\langle \text{root}, q, S_1 \rangle \wedge \text{root} \neq q \wedge S = \{v\} \sqcup S_1 \wedge (\forall x \in S_1. v \leq x))$$

where the parameters  $p$  and  $q$  denote the `prev` field of `root` and the `next` of the list's last node, respectively. Meanwhile  $S$  is a bag (multi-set) parameter to represent the list's content. We can see in the base case of definition that  $S = \emptyset$ , and in the recursive case that all values stored after `root` must be no less than `root`'s value.

Another example is the definition of node-balanced trees with binary search property:

$$\begin{aligned} \text{nbt}\langle S \rangle \equiv & (\text{root} = \text{null} \wedge S = \emptyset) \vee \\ & (\text{root}::\text{node2}\langle v, p, q \rangle * p::\text{nbt}\langle S_p \rangle * q::\text{nbt}\langle S_q \rangle \wedge S = \{s\} \sqcup S_p \sqcup S_q \wedge \\ & (\forall x \in S_p. x \leq s) \wedge (\forall x \in S_q. s \leq x) \wedge -1 \leq |S_p| - |S_q| \leq 1) \end{aligned}$$

where  $S$  captures the content of the tree. We require the difference in node numbers of the left and right sub-trees be within one, as the node-balanced property indicates.

User-defined predicates may then be used to specify loop invariants and method pre/post-specifications. In HIP/SLEEK, the HIP verifier is used to automatically verify programs against their specifications, while the SLEEK prover is invoked by the verifier to conduct entailment proofs. Given two separation formulas  $\Delta_1$  and  $\Delta_2$ , SLEEK attempts to prove that  $\Delta_1$  entails  $\Delta_2$ ; if it succeeds, it returns a frame  $R$  such that  $\Delta_1 \vdash \Delta_2 * R$ . For instance, given the entailment

$$p::ll\langle n \rangle \wedge n > 0 \vdash \exists q. p::node\langle q \rangle$$

SLEEK produces the following result after unfolding the LHS predicate:

$$p::ll\langle n \rangle \wedge n > 0 \vdash \exists q. p::node\langle q \rangle * [q::ll\langle n-1 \rangle \wedge n > 0]$$

where the inferred frame, or residue, is shown in squared brackets. The proposed analysis in this paper will use SLEEK to perform deductions of separation formulas.

## 2.2 An Illustrative Example

We illustrate our approach using method `insert_sort` in Fig 1. We show how our analysis infers missing constraints to improve the user-supplied incomplete specification, and how it analyses auxiliary methods without user-annotations.

<pre> 1 data node { int val; node next; } 2 node insert_sort(node x) 3   requires x::llB(S) 4   ensures res::sllB(T) { 5     if (x.next == null) return x; 6     else { node s = x.next; 7       node r = insert_sort(s); 8       return insert(r, x); 9     } 10  }</pre>	<pre> 11 node insert(node r, node x) { 12   if (r == null) { 13     x.next = null; return x; 14   } else if (x.val &lt;= r.val) { 15     x.next = r; return x; 16   } else { 17     r.next = insert(r.next, x); 18     return r; 19   } 20 }</pre>
--	--

**Fig. 1.** The insertion sort program for lists.

The `insert_sort` method sorts a singly-linked list. It takes in an unsorted list starting from  $x$  with content  $S$  and returns a sorted list (lines 3 and 4 where `res` denotes the method return value). The algorithm first sorts the list referenced by `x.next` recursively (line 7), and then inserts node  $x$  into the resulted sorted list (line 8). For the node insertion, it invokes another method `insert` for which

the user has not provided a specification. We call `insert` an auxiliary method and `insert_sort` a primary one.

For the primary method with a partial specification, our analysis proceeds in two steps. Firstly, starting from the partial precondition, a forward analysis is conducted to compute the postcondition of the method in the form of a *constraint abstraction* [21]. This constraint abstraction is effectively a transfer function for the method, which may be recursively defined. During this analysis, abductive reasoning may be used whenever the current state fails to establish the precondition of the next program command. Secondly, instead of a direct fixpoint computation in the combined abstract domain (with shape, numerical and bag information), a “pure” constraint abstraction (without heap shape information) is derived from the generated constraint abstraction and the user-given partial postcondition. This pure constraint abstraction is then solved by fixpoint solvers in pure (numerical/set) domains, such as [36, 38].

The constraint abstraction of a code segment (e.g. a method) in our settings is an abstraction form of that code’s postcondition, given a certain precondition. As the code may contain loops or recursive calls, its constraint abstraction can also be recursive, or in an *open form*, accordingly. To illustrate, for the following while loop and its precondition

$$\{x \geq 0 \wedge y = 0\} \text{ while } (x > 0) \{x = x - 1; y = y + 1;\}$$

we have its constraint abstraction as

$$Q(x, x', y, y') ::= x = 0 \wedge x = x' \wedge y = y' \vee x > 0 \wedge Q(x - 1, x', y + 1, y')$$

where we denote  $x$  and  $y$  as their values before the loop, and the primed versions as the values after the loop execution (we will explain this in more detail in Sec 3). Such constraint abstraction presents the postcondition of the while loop. Its fixpoint can be achieved with a standard fixpoint calculation process, with result  $x \geq 0 \wedge y = 0 \wedge x' = 0 \wedge y' = x$ . However, as will be seen later, our constraint abstraction is generally more complicated involving both shape and pure constraints, requiring us to split them for solution somehow.

As for the example, our forward analysis runs on the body of `insert_sort` to construct the constraint abstraction. For lines 5-9, it produces a disjunction as the effect of if-else (according to the if-else rule in page 30):

$$Q(x, S, \text{res}, T) ::= (\text{post-state of if}) \vee (\text{post-state of else})$$

where  $Q$  represents the post-state of the if-else statement (as well as the method), and its parameters  $x, S, \text{res}$  and  $T$  are the (program and logical) variables involved in the state.

For the if branch, after the unfolding over  $x::11B(S)$  (rule `unfold` in page 28), we know from the condition that the input list  $x$  has only one node, and thus its post-state will be

$$\exists v \cdot x::\text{node}(v, \text{null}) \wedge \text{res} = x \wedge S = \{v\} \quad (1)$$

Meanwhile, for the else branch, the list will firstly be unrolled by one node at line 6 (rule `unfold`), making `x.next` point to `s` (rule `assign` in page 30), which references a sub-list one node shorter than the input list beginning from `x`:

$$\exists S_s, v \cdot x::\text{node}\langle v, s \rangle * s::\text{llb}\langle S_s \rangle \wedge S = S_s \sqcup \{v\} \quad (2)$$

After that, `insert_sort` is invoked recursively with `s`. It will consume the precondition (`s::llb⟨Ss⟩`) and ensure the postcondition (in terms of `Q`, partially according to the rule in page 29; however it will be substituted as described later). In that case, the state immediately after symbolic execution of line 7 is

$$\begin{aligned} Q(x, S, \text{res}, T) ::= & \exists v \cdot x::\text{node}\langle v, \text{null} \rangle \wedge \text{res} = x \wedge S = \{v\} \vee \\ & \exists v, s, S_s, r, S_r \cdot x::\text{node}\langle v, s \rangle * Q(s, S_s, r, S_r) \wedge |S| > 1 \wedge S = S_s \sqcup \{v\} \end{aligned}$$

Note that existential variables (not in the parameter list of `Q`) are local variables whose quantification may be omitted for brevity. The first disjunctive branch corresponds to the base case in the method body, and the second branch captures the effect of the recursive call (with `Q`).

Then the forward analysis continues over line 8 to invoke `insert`. Because the user has provided no annotations for that method, its specifications must be synthesised. For this purpose we replace `Q(s, Ss, r, Sr)` in second branch with `r::sllb⟨Sr⟩ ∧ P(s, Ss, r, Sr)` to make explicit the heap portion referred to by `r` before we analyse the auxiliary call `insert(r, x)` (rule `call-inf` in page 29). This is safe because the following entailment relationship is added to our assumption:

$$Q(x, S, \text{res}, T) \vdash \text{res}::\text{sllb}\langle T \rangle \wedge P(x, S, \text{res}, T) \quad (3)$$

which signifies that `Q` can be abstracted as a sorted list referenced by `res` plus some pure constraints `P` (also in constraint abstraction form, whose definition is to be derived in the next step). and hence `insert`'s precondition can be figured out from the symbolic state at call site, and its postcondition will be computed as well. The analysis for such auxiliary methods (including loops) works in the same way as that for primary methods, except that a pre-analysis is involved to figure out the raw pre/post shape information (before invoking the analysis algorithm for primary methods). More details are explained slightly later.

$$\text{requires } r::\text{sllb}\langle S \rangle * x::\text{node}\langle v, \_ \rangle \text{ ensures } \text{res}::\text{sllb}\langle T \rangle \wedge T = S \sqcup \{v\} \quad (4)$$

which indicates that the returned list has the same content as the input list (`x`) plus `{v}`. Applying it, we obtain the following post-state for `insert_sort`:

$$\begin{aligned} Q(x, S, \text{res}, T) ::= & x::\text{node}\langle v, \text{null} \rangle \wedge \text{res} = x \wedge S = \{v\} \vee \\ & \text{res}::\text{sllb}\langle S_{\text{res}} \rangle \wedge P(s, S_s, r, S_r) \wedge |S| > 1 \wedge S = S_s \sqcup \{v\} \wedge S_{\text{res}} = S_r \sqcup \{v\} \end{aligned}$$

The first disjunctive branch corresponds to the base case, but the second branch now captures the effect of the recursive call as well as the auxiliary call (to `insert`). In the base case, the method's return pointer (`res`) points to one node with value `v`. The recursive branch signifies that the post-state of the method

concerns the recursive call and the auxiliary call (over  $\mathbf{s}$  and  $\mathbf{r}$ ), as the constraint abstraction denotes. Note that  $\mathbf{T}$  will be not available (as well as its relationship with  $\mathbf{S}_{\text{res}}$ ) until next step.

In the second step, we first derive the definition of the pure constraint abstraction  $\mathbf{P}$  from the above post-state  $\mathbf{Q}$ . Each disjunctive branch of  $\mathbf{Q}$  is used to entail the user-given post-shape (with appropriate instantiations of the parameters). The obtained frames form (via disjunction) the definition of  $\mathbf{P}$ . For `insert.sort`, we obtain the following pure constraint abstraction:

$$\mathbf{P}(\mathbf{x}, \mathbf{S}, \text{res}, \mathbf{T}) ::= (\mathbf{T}=\mathbf{S} \wedge |\mathbf{S}|=1) \vee (\mathbf{P}(\mathbf{s}, \mathbf{S}_s, \mathbf{r}, \mathbf{S}_r) \wedge |\mathbf{S}|>1 \wedge \mathbf{S}=\mathbf{S}_s \sqcup \{\mathbf{v}\} \wedge \mathbf{T}=\mathbf{S}_r \sqcup \{\mathbf{v}\})$$

We then use pure fixpoint solvers to obtain a closed-form formula  $|\mathbf{S}|\geq 1 \wedge \mathbf{T}=\mathbf{S}$  for  $\mathbf{P}$ . Based on (3), we now obtain the closed-form approximation for  $\mathbf{Q}$ :

$$\mathbf{Q}(\mathbf{x}, \mathbf{S}, \text{res}, \mathbf{T}) ::= \text{res}::\text{s11B}(\mathbf{T}) \wedge |\mathbf{S}|\geq 1 \wedge \mathbf{T}=\mathbf{S}$$

The obtained pure formula is then used to refine the method's specification as

$$\text{requires } \mathbf{x}::\text{l1B}(\mathbf{S}) \wedge |\mathbf{S}|\geq 1 \quad \text{ensures } \text{res}::\text{s11B}(\mathbf{T}) \wedge \mathbf{T}=\mathbf{S}$$

which imposes more requirement in the precondition, stating that there should be at least one node in the list to be sorted for the sake of memory safety. With that obligation, the method guarantees that the result list is sorted and its content remains the same as the input list.

### 2.3 Analysis for the Unannotated Method in Example

The unannotated method `insert` in the example inserts a node  $\mathbf{x}$  into a sorted list  $\mathbf{r}$ . It judges three cases and has a non-tail-recursive call to itself in the last case (to insert  $\mathbf{x}$  after list  $\mathbf{r}$ 's head). As we want to minimise user's annotations, we do not require the user to supply loop invariant; instead we will calculate the loop's postcondition. Since no user-annotations are provided, our analysis synthesises its (raw) pre- and post-shapes which are then refined in the same way as for primary methods. The pre-shape is directly synthesised from the abstract program state at the call site ( $\mathbf{x}::\text{node}(\mathbf{v}, \mathbf{s}) * \mathbf{r}::\text{s11B}(\mathbf{S}_r)$ ). We unroll the recursive call once, symbolically execute the unrolled method body (starting from the pre-shape) to obtain a post-state, and then use the post-state to filter out any invalid post-shapes from the set of possible post-shapes (drawn from all available shape predicates). For this example, the possible post-shapes can be (a)  $\mathbf{x}::\text{s11B}(\mathbf{S}_1) * \text{res}::\text{s11B}(\mathbf{S}_2)$ , and (b)  $\text{res}::\text{s11B}(\mathbf{S})$ , etc. The symbolic execution gives the following post-state:

$$\begin{aligned} & \mathbf{x}::\text{node}(\mathbf{v}, \text{null}) \wedge \mathbf{x}=\text{res} \vee \mathbf{x}::\text{node}(\mathbf{v}, \mathbf{r}) * \mathbf{r}::\text{s11B}(\mathbf{S}_1) \wedge \mathbf{x}=\text{res} \wedge (\forall \mathbf{u} \in \mathbf{S}_1. \mathbf{v} \leq \mathbf{u}) \vee \\ & \quad \mathbf{r}::\text{node}(\mathbf{u}, \mathbf{x}) * \mathbf{x}::\text{node}(\mathbf{v}, \text{null}) \wedge \mathbf{r}=\text{res} \wedge \mathbf{u} \leq \mathbf{v} \vee \\ & \quad \mathbf{r}::\text{node}(\mathbf{u}, \mathbf{x}) * \mathbf{x}::\text{node}(\mathbf{v}, \mathbf{r}_1) * \mathbf{r}_1::\text{s11B}(\mathbf{S}_1) \wedge \mathbf{r}=\text{res} \wedge \mathbf{u} \leq \mathbf{v} \wedge (\forall \mathbf{w} \in \mathbf{S}_1. \mathbf{v} \leq \mathbf{w}) \end{aligned}$$

which does not entail the candidate (a), so we filter it out. Taking (b) as the post-shape, we can employ the same analysis for the primary method to obtain the specification (4) (page 7) for `insert` and continue with the analysis for the primary method.



## 2.4 Another Illustrative Example

We illustrate our approach with another more interesting example. We show how the user is expected to provide shape information for specifications of a primary method, and how our proposed analysis will refine such specifications with pure constraints, and derive specifications for loops without annotations.

```

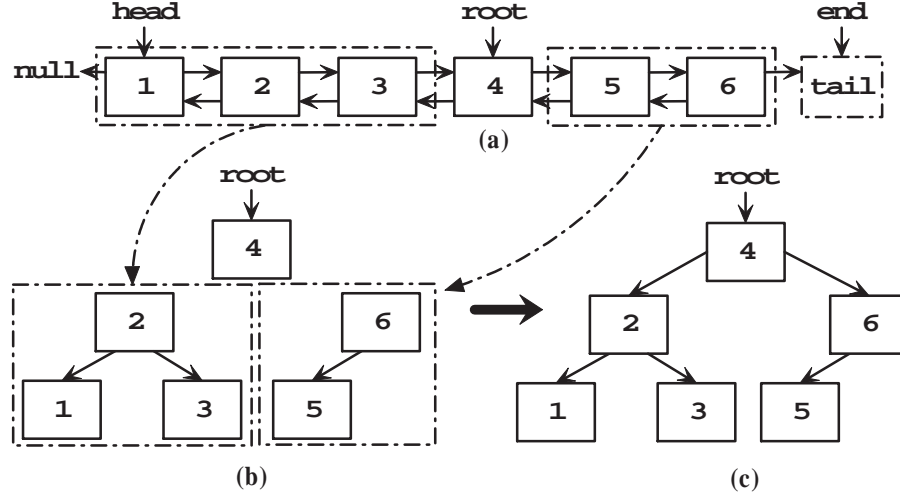
0 data node2 { int val;           14 if (head == root)
    node2 prev; node2 next; }    15     root.prev = null;
1 node2 sdl2nbt(node2 head,      16 else
    node2 tail)                 17     root.prev = sdl2nbt(head, root);
2 requires head::sdlB(p,q,S)     18 node2 tmp = root.next;
3 ensures res::nbt(Sres)         19 if (tmp == tail)
4 {                               20     root.next = null;
5 node2 root = head;            21 else {
6 node2 end = head;             22     tmp.prev = null;
7 while(end != tail) {          23     root.next = sdl2nbt(tmp, tail);
8     end = end.next;           24 }
9     if (end != tail) {        25 return root;
10         end = end.next;       26 }
11         root = root.next;
12     }
13 }

```

**Fig. 2.** The sorted doubly-linked list to node-balanced tree method.

Let us consider the method `sdl2nbt` shown in Fig 2. Taking in a sorted doubly-linked list as argument, `sdl2nbt` will convert it into a node-balanced tree together with binary search properties, as indicated in lines 2 and 3. Its algorithm proceeds as follows: first it finds the “centre” node in the list (`root`), where the difference of numbers of its left and right nodes is at most one, as Fig 3 (a) indicates (lines 5-13). Then it applies the algorithm recursively on both list segments on the centre’s left and right hand sides, and regards the centre node as the tree’s root, whose left and right children are the resulted subtrees’ roots from the recursive calls, as in Fig 3 (b) and (c) (lines 14-25). As the data structure of doubly-linked list and binary tree are homomorphic (line 0), we reuse the nodes in the input list instead of creating a new tree, making this algorithm in-place. The parameter `head` in line 1 denotes the first node of the input list, and `tail` is where the list’s last node’s `next` field points to. When using this method `tail` should be set as `null` initially.

Our framework allows the user to verify and/or refine a number of properties about this code. Firstly, the transformation of shapes from initial to final states (namely, from a doubly-linked list to a binary tree) must be captured. Secondly, some structural numerical information should be inferred, so as to prove the node counts before and after the method invocation are the same and the node-



**Fig. 3.** Transferring from a sorted doubly-linked list to a node-balanced BST.

balanced property of the tree, etc. Meanwhile, we also want to derive relational numerical information as lists' sortedness and trees' binary search property, and finally set/bag information like the symbolic content of the list's and the tree's (in order to prove the values stored in the list and the resulted tree are the same). Finally, some obligation for memory safety should be found in the precondition, to ensure the input list is non-empty (otherwise the dereference in line 15/17 will fail). To deal with all these properties, we expect the user to provide shape information for primary methods' specifications as in Fig 2. Based on that, we try to compute the remaining constraints (including the missing parts of pure specifications for primary methods, and both pre- and post-conditions for loops and auxiliary methods).

As for the example, as the user has provided the pre- and post-shapes for method `sd12nbt`, our analysis proceeds in two steps: generating the constraint abstraction, and solving it. The first step is mainly a symbolic execution over the program to find its postcondition, so as to generate the constraint abstraction. During this step, for any loops and/or auxiliary method calls (lines 7-13 in the example), the symbolic execution will invoke the analysis procedure again to compute their specifications for the current execution to continue (Section 2.5). Therefore we can find the while loop's postcondition as

$$\begin{aligned} & \text{head}::\text{sd1B}(\text{null}, \text{root}, S_h) * \text{root}::\text{sd1B}(p, \text{tail}, S_r) \wedge \\ & \text{end}=\text{tail} \wedge S=S_h \sqcup S_r \wedge (\forall x \in S_h, y \in S_r. x \leq y) \wedge \underline{0 \leq |S_h| - |S_r| \leq 1} \end{aligned} \quad (5)$$

which indicates that the original list segment starting from **head** is cut into two pieces with a cutpoint **root**, where both are still sorted and the content is also preserved. Meanwhile, the essential constraint (the underlined part, saying the

list beginning with `head` is at most one node longer than that with `root`) to ensure the node-balanced property is derived as well.

When the symbolic execution finishes, it generates the following constraint abstraction as the postcondition of the method:

$$\begin{aligned}
Q(\text{head}, p, q, S, \text{res}, S_{\text{res}}) ::= & \\
& \text{root}::\text{node2}\langle v, \text{null}, \text{null} \rangle \wedge \text{head}=\text{root}=\text{res} \wedge \text{tmp}=\text{q}=\text{tail} \wedge p=\text{null} \wedge \\
& S=\{v\} \vee \\
& \text{head}::\text{node2}\langle s, \text{null}, \text{root} \rangle * \text{root}::\text{node2}\langle v, \text{res}_h, \text{null} \rangle \wedge \text{res}=\text{root} \wedge \\
& \text{tmp}=\text{q}=\text{tail} \wedge p=\text{null} \wedge S=\{s, v\} \wedge s \leq v \vee \\
& \text{res}_h::\text{nbt}\langle S_{\text{res}}^h \rangle * \text{res}_r::\text{nbt}\langle S_{\text{res}}^r \rangle * \text{root}::\text{node2}\langle v, \text{res}_h, \text{res}_r \rangle \wedge \\
& P(\text{head}, p, \text{root}, S_h, \text{res}_h, S_{\text{res}}^h) \wedge P(\text{tmp}, \text{null}, \text{tail}, S_r, \text{res}_r, S_{\text{res}}^r) \wedge \\
& \text{head} \neq \text{root} \wedge \text{root}=\text{res} \wedge \text{tmp} \neq \text{tail} \wedge q=\text{tail} \wedge \\
& S=S_h \sqcup \{v\} \sqcup S_r \wedge (\forall x \in S_h, y \in S_r \cdot x \leq v \leq y) \wedge 0 \leq |S_h| - |S_r| \leq 1
\end{aligned}$$

where  $P$  stands for corresponding pure constraint abstraction explained below. The first two disjunctive branches are base cases of the method's invocation, and the last denotes the effect of recursive calls combined into the postcondition. The first case represents the scenario where there is only one node in the original list (with `res` as the method's return value). The second is for the case of two nodes, one referenced by `head`, pointing to the other one, `root`. In this case the value of `head` is no more than that of `root`. The third case is defined recursively with the constraint abstraction itself, meaning that the post-state concerns the `root` node and the post-states of two recursive calls over `head` and `tmp`, respectively. Note that  $S_{\text{res}}$  does not appear in  $Q$ 's definition. Since it stands for pure properties in user-provided post-shape, it will be involved when we abstract  $Q$  against that post-shape in the next step.

The second step solves the constraint abstraction  $Q$  by finding a closed-form approximation of it. Instead of performing a fixpoint analysis directly on  $Q$  over the combined domain, we first derive a pure constraint abstraction  $P$  (with the help of SLEEK) from  $Q$  and the user-provided heap part of postcondition. Then we are able to use some existing conventional solvers [36, 38] to compute the pure fixpoint. For the `sdl2nbt` method, we generate the pure constraint abstraction  $P$  based on the following entailment relation:

$$Q(\text{head}, p, q, S, \text{res}, S_{\text{res}}) \vdash \text{res}::\text{nbt}\langle S_{\text{res}} \rangle \wedge P(\text{head}, p, q, S, \text{res}, S_{\text{res}})$$

which produces the following pure constraint abstraction  $P$ :

$$\begin{aligned}
P(\text{head}, p, q, S, \text{res}, S_{\text{res}}) ::= & \\
& \text{head}=\text{root}=\text{res} \wedge \text{tmp}=\text{q}=\text{tail} \wedge p=\text{null} \wedge S=S_{\text{res}}=\{v\} \vee \\
& \text{head} \neq \text{root} \wedge \text{res}=\text{root} \wedge \text{tmp}=\text{q}=\text{tail} \wedge p=\text{null} \wedge \\
& S=S_{\text{res}}=\{s, v\} \wedge s \leq v \vee \\
& P(\text{head}, p, \text{root}, S_h, \text{res}_h, S_{\text{res}}^h) \wedge P(\text{tmp}, \text{null}, \text{tail}, S_r, \text{res}_r, S_{\text{res}}^r) \wedge \\
& \text{head} \neq \text{root} \wedge \text{root}=\text{res} \wedge \text{tmp} \neq \text{tail} \wedge q=\text{tail} \wedge S=S_h \sqcup \{v\} \sqcup S_r \wedge \\
& S_{\text{res}}=S_{\text{res}}^h \sqcup \{v\} \sqcup S_{\text{res}}^r \wedge (\forall x \in S_h, y \in S_r \cdot x \leq v \leq y) \wedge 0 \leq |S_h| - |S_r| \leq 1
\end{aligned}$$

Note that the heap information is already eliminated from  $P$ ; instead the constraints over  $S_{\text{res}}$  are included during the entailment checking procedure. This allows us to solve  $P$  to refine the user-provided shape-only specification.

After solving  $P$ , we achieve the following constraint:

$$p=\text{null} \wedge q=\text{tail} \wedge S=S_{\text{res}} \wedge |S| \geq 1$$

with which we can refine the method's specifications as

$$\begin{aligned} \text{requires } & \text{head}::\text{sdlb}\langle p, q, S \rangle \wedge p=\text{null} \wedge q=\text{tail} \wedge |S| \geq 1 \\ \text{ensures } & \text{res}::\text{nbt}\langle S_{\text{res}} \rangle \wedge S=S_{\text{res}} \end{aligned}$$

which proposes more requirements in the precondition, as the `head`'s `prev` field should be `null`, and the whole list's last node's `next` field must point to `tail`. Meanwhile, there should be at least one node in the list for the sake of memory safety. With those obligations, the method guarantees that the result is an node-balanced tree with binary search property, whose content is the same as the input list.

## 2.5 Analysis for the While Loop in Example

In the method `sd12nbt`, there is a while loop (lines 7-13) to discover the centre node of the given list segment referenced by `head`. It traverses the list segment with two pointers `root` and `end`. The `end` pointer goes towards the list segment's tail twice as fast as `root`. When `end` arrives at the tail of the segment (`tail`), `root` will point to the list segment's centre node.

As we want to minimise user's annotations, we do not require the user to supply loop invariant; instead we will try to calculate its postcondition. As aforementioned, our analysis must first synthesise its pre- and post-states with shape information, and then proceed with its constraint abstraction. For pre it is straightforward as the program state before the loop will provide relevant shape information. For post it is done by checking the loop body (unrolled once)'s symbolic execution result against all possible abstracted shapes. For the previous example, we first generate all possible shapes according to the variables accessed by the loop, such as (a) `head::sdlb` $\langle p_h, q_h, S_h \rangle * \text{root}::\text{sdlb}\langle p_r, q_r, S_r \rangle$ , and (b) `head::sdlb` $\langle p_h, q_h, S_h \rangle * \text{root}::\text{nbt}\langle h_r, b_r, S_r \rangle$ , and many so forth. Then the unrolled loop body is symbolically executed several times to filter out any invalid shape as an invariant. In the example's case, executing the loop body will yield the following result:

$$\begin{aligned} & \text{head}::\text{node2}\langle v, p, \text{end} \rangle \wedge \text{head}=\text{root} \wedge \text{end}=\text{tail} \vee \\ & \text{head}::\text{node2}\langle v_h, p, \text{root} \rangle * \text{root}::\text{node2}\langle v_r, \text{head}, \text{end} \rangle \wedge \text{end}=\text{tail} \end{aligned} \quad (6)$$

where (b) is directly filtered out since  $(6) \vdash (b) * \text{true}$  fails. However (a) remains a candidate, as both  $(6) \vdash (a) * \text{true}$  holds. Therefore, regarding (a) as a possible shape post, we can employ the same approach for the whole method to generate

a constraint abstraction for the while loop, and solve it to achieve formula (6) in the last section.

One more note for the while loop in this example is that the symbolic execution may actually permit more than one shapes to enter as candidates, e.g.  $\text{head}::\text{sd1B}\langle p_h, q_h, S_h \rangle$ . Generally this does not affect the analysis result, as we allow the analysis to continue with all possible postconditions computed from this while loop, and always choose the most precise final result. In the motivating example, both  $\text{head}::\text{sd1B}\langle p_h, q_h, S_h \rangle$  and (a) are valid shape postconditions for the loop, but later the former one will cause the analysis to fail in line 15/17, because it inappropriately approximated the invariant and hence lost information about  $\text{root}$ . Since we synthesise all possible shapes, we can always select those shapes sufficiently strong to support further analysis to obtain a meaningful result.

### 3 Language and Abstract Domain

To simplify presentation, we focus on a strongly-typed C-like imperative language in Fig 4. A program  $\text{Prog}$  consists of type declarations  $tdecl$ , which can define either data type  $\text{datatype}$  (e.g.  $\text{node}$ ) or predicate  $\text{spred}$  (e.g.  $\text{11B}$ ), and some method declarations  $\text{meth}$ . The definitions for  $\text{spred}$  and  $\text{mspec}$  are given later in Fig 5. The language is expression-oriented, so the body of a method is an ex-

$\text{Prog} ::= tdecl^* \text{meth}^*$	$tdecl ::= \text{datatype} \mid \text{spred} \mid \text{lemma}$
$\text{datatype} ::= \text{data } c \{ \text{field}^* \}$	$\text{field} ::= t \ v \quad t ::= c \mid \tau$
$\text{meth} ::= t \ mn \ ((t \ v)^*; (t \ v)^*) \ \text{mspec}^* \{e\}$	$\tau ::= \text{int} \mid \text{bool} \mid \text{void}$
$e ::= d \mid d[v] \mid v=e \mid e_1; e_2 \mid t \ v; \ e \mid \text{if } (v) \ e_1 \ \text{else } e_2 \mid \text{while } (v) \{e\}$	
$d ::= \text{null} \mid k^\tau \mid v \mid \text{new } c(v^*) \mid mn(u^*; v^*)$	
$d[v] ::= v.f \mid v.f:=w \mid \text{free}(v)$	

**Fig. 4.** A Core (C-like) Imperative Language.

pression composed of  $e$  (the recursively defined program constructor) and  $d$  and  $d[v]$  (atom instructions). We also allow both call-by-value and call-by-reference method parameters (which are separated with a semicolon ; where the ones before ; are call-by-value and the ones after are call-by-reference).

Our specification language (in Fig 5) allows (user-defined) shape predicates to specify both separation and pure properties. The shape predicates  $\text{spred}$  and lemmas  $\text{lemma}$  are constructed with disjunctive constraints  $\Phi$ . We require that the predicates be well-formed [35]. A conjunctive abstract program state,  $\sigma$ , is composed of a heap (shape) part  $\kappa$  and a pure part  $\pi$ , where  $\pi$  consists of  $\gamma, \phi$  and  $\varphi$  as aliasing, numerical and bag information, respectively. We use  $\text{SH}$  to denote the set of such conjunctive states. During the symbolic execution, the abstract program state at each program point will be a disjunction of  $\sigma$ 's, denoted by  $\Delta$ . Note that constraint abstractions (e.g.  $\mathbb{Q}(v^*)$ ) may occur in  $\Delta$  during the analysis.

$spread$	$::= pred(v^*) \equiv \Phi$	$lemma$	$::= pred(v^*) \wedge \pi \leftarrow \Phi$
$mspec$	$::= requires \Phi_{pr} ensures \Phi_{po}$		
$\Delta$	$::= Q(v^*) \mid \Phi \mid \Delta_1 \vee \Delta_2 \mid \Delta \wedge \pi \mid \Delta_1 * \Delta_2 \mid \exists v. \Delta$		
$\Phi$	$::= \bigvee \sigma^* \mid \sigma ::= \exists v^*. \kappa \wedge \pi$		
$\Upsilon$	$::= P(v^*) \mid \bigvee \omega^* \mid \Upsilon_1 \wedge \Upsilon_2 \mid \Upsilon_1 \vee \Upsilon_2 \mid \exists v. \Upsilon$		
$\kappa$	$::= \mathbf{emp} \mid v \mapsto c(v^*) \mid pred(v^*) \mid \kappa_1 * \kappa_2$		
$\omega$	$::= \exists v^*. \pi \mid \pi ::= \gamma \wedge \phi$		
$\gamma$	$::= v_1 = v_2 \mid v = \mathbf{null} \mid v_1 \neq v_2 \mid v \neq \mathbf{null} \mid \gamma_1 \wedge \gamma_2$		
$\phi$	$::= \varphi \mid b \mid a \mid \phi_1 \wedge \phi_2 \mid \phi_1 \vee \phi_2 \mid \neg \phi \mid \exists v. \phi \mid \forall v. \phi$		
$b$	$::= \mathbf{true} \mid \mathbf{false} \mid v \mid b_1 = b_2 \mid a ::= s_1 = s_2 \mid s_1 \leq s_2$		
$s$	$::= k^{\mathbf{int}} \mid v \mid k^{\mathbf{int}} \times s \mid s_1 + s_2 \mid -s \mid \max(s_1, s_2) \mid \min(s_1, s_2) \mid  B $		
$\varphi$	$::= v \in B \mid B_1 = B_2 \mid B_1 \sqsubseteq B_2 \mid B_1 \sqsubset B_2 \mid \forall v \in B. \phi \mid \exists v \in B. \phi$		
$B$	$::= B_1 \sqcup B_2 \mid B_1 \sqcap B_2 \mid B_1 - B_2 \mid \{\} \mid \{v\}$		

Fig. 5. The Specification Language.

A closed-form  $\Delta$  (containing no constraint abstractions) can be normalised to the  $\Phi$  form [35]. Pure constraint abstraction  $P$  is analogously defined to  $Q$ .

The memory model of our specification formulae is adapted from the model given for “early versions” of separation logic [40], except that we have extensions to handle user-defined shape predicates and related pure properties. We assume sets  $\mathbf{Loc}$  of memory locations,  $\mathbf{Val}$  of primitive values (with  $0 \in \mathbf{Val}$  denoting  $\mathbf{null}$ ),  $\mathbf{Var}$  of variables (program and logical variables), and  $\mathbf{ObjVal}$  of object values stored in the heap, with  $c[f_1 \mapsto \nu_1, \dots, f_n \mapsto \nu_n]$  denoting an object value of data type  $c$  where  $\nu_1, \dots, \nu_n$  are current values of the corresponding fields  $f_1, \dots, f_n$ . Let  $s, h \models \Delta$  denote the model relation, i.e. the stack  $s$  and heap  $h$  satisfy  $\Delta$ , with  $h, s$  from the following concrete domains:

$$h \in \mathbf{Heaps} =_{df} \mathbf{Loc} \rightarrow_{fin} \mathbf{ObjVal} \quad s \in \mathbf{Stacks} =_{df} \mathbf{Var} \rightarrow \mathbf{Val} \cup \mathbf{Loc}$$

Note that each heap  $h$  is a finite partial mapping while each stack  $s$  is a total mapping, as in the classical separation logic [24, 40]. The detailed model definitions can be found in Nguyen et al. [35].

In the analysis we use three kinds of variables in the  $\mathbf{Var}$  set: program variables, logical variables related to program variables’ shapes (such as a list’s length), and logical variables to record intermediate states. For the first two groups we use variables without subscription (such as  $\mathbf{x}$  and  $\mathbf{xn}$ ), and denote a program variable’s initial value as unprimed, and its current (and hence final) value as primed [9, 35]. For the third group, we use subscript ones like  $\mathbf{x}_1$  and  $\mathbf{xn}_1$ . For instance, for a code segment  $\mathbf{x} := \mathbf{x} + 1; \mathbf{x} := \mathbf{x} - 2$  starting with state  $\{\mathbf{x} > 1\}$ , we have the following reasoning procedure:

$$\{\mathbf{x}' = \mathbf{x} \wedge \mathbf{x} > 1\} \mathbf{x} := \mathbf{x} + 1 \{\mathbf{x} > 1 \wedge \mathbf{x}' = \mathbf{x} + 1\} \mathbf{x} := \mathbf{x} - 2 \{\mathbf{x} > 1 \wedge \mathbf{x}' = \mathbf{x}_1 - 2 \wedge \mathbf{x}_1 = \mathbf{x} + 1\}$$

where the final value of  $\mathbf{x}$  is recorded in variable  $\mathbf{x}'$  and  $\mathbf{x}_1$  keeps an intermediate state of  $\mathbf{x}$ .

## 4 The Analysis

The overall algorithm is listed in Fig 6.

```

Algorithm Analysis( $\mathcal{T}, \mathcal{S}, mn, \sigma, x^*, y^*$ )
1  case  $mn$  of
2    | while  $(w) \{e_0\} \rightarrow f := \text{fresh\_name}(); e := \text{if } (w) \{e_0; f(x^*; y^*)\};$ 
       $(u^*, v^*) := (x^*, y^*); ([(\Phi_{pr}^i, \Phi_{po}^i)], n) := \text{Preproc}(\mathcal{T}, \mathcal{S}, f, x^*, y^*, e_0, \sigma, x^*, y^*);$ 
       $prim := \text{false};$ 
3    |  $t \ mn \ ((t \ u_0)^*; (t \ v_0)^*) \{e_0\} \rightarrow f := mn; e := e_0; (u^*, v^*) := (u_0^*, v_0^*);$ 
       $([(\Phi_{pr}^i, \Phi_{po}^i)], n) := \text{Preproc}(\mathcal{T}, \mathcal{S}, f, u^*, v^*, e_0, \sigma, x^*, y^*); prim := \text{false};$ 
4    |  $t \ mn \ ((t \ u_0)^*; (t \ v_0)^*) \ (\text{requires } \Phi_{pr}^i \ \text{ensures } \Phi_{po}^i)_{i=1}^m \{e_0\} \rightarrow f := mn;$ 
       $e := e_0; (u^*, v^*) := (u_0^*, v_0^*); n := m; (\Phi_{pr}^i, \Phi_{po}^i)_{i=1}^m := (\Phi_{pr}^i, \Phi_{po}^i)_{i=1}^m;$ 
       $prim := \text{true};$ 
5  end case
6   $sps := \emptyset$ 
7  for  $i := 1$  to  $n$  do
8     $sp := \text{CA\_Gen\_Solve}(\mathcal{T}, f, e, \Phi_{pr}^i, \Phi_{po}^i, u^*, v^*)$ 
9    if  $prim = \text{false}$  and  $sp \neq \text{fail}$  then return  $(f, sp)$ 
10   else if  $prim = \text{true}$  then  $sps := sps \cup sp$ 
11   end if
12 end for
13 return  $(f, sps)$ 
end Algorithm

```

**Fig. 6.** Main analysis algorithm.

Our analysis algorithm takes as input all available specifications and shapes, and the code segment to be analysed, together with an optional conjunctive program state and two variable sequences (mainly for loops and auxiliary methods). The algorithm first recognises the type of input code segment ( $mn$  in line 1). In the first two cases (while loop in line 2 and auxiliary method call in line 3), we do not know anything about the code's specifications; therefore some preprocessing should be done to discover the code's pre- and post-shapes with **Preproc** (Fig. 7). For primary method (line 4), as the user should have provided the shape-based specifications, no preprocessing is needed. Then the constraint abstraction generation and solving algorithm is applied to each specification to refine it (line 8). Note here we apply a lazy scheme for loops and auxiliary methods: as the pre-processing may yield several possible shape specifications in a list (ordered with heuristics such that the specifications with more possibility to make the whole verification succeed are more in front), we try to verify each in sequence. Once a specification can be verified against the program, then it is returned and the other ones are omitted. In this way we try to make our verification more scalable, as still will be described in later sections.

```

Algorithm Preproc( $\mathcal{T}, \mathcal{S}, f, u^*, v^*, e, \sigma, x^*, y^*$ )
1   $sps := []$ ;
2   $prs := \text{SynPre}(\mathcal{S}, f, u^*, v^*, \sigma, x^*, y^*)$ 
3  for  $\Phi_{pr} \in prs$  do
4     $pos := \text{SynPost}(\mathcal{T}, \mathcal{S}, f, e, \Phi_{pr}, u^*, v^*)$ 
5     $sps := \text{concat}(sps, pos)$ 
6  end for
7  return  $(sps, |sps|)$ 
end Algorithm

```

**Fig. 7.** Pre-processing algorithm.

The pre-processing algorithm mainly invokes the shape synthesis procedures to discover all possible pre- and post-shapes for loops and auxiliary methods, as shown in lines 1 and 4. Then the list of shape pairs (specifications) are returned and used in further analysis. The details of shape synthesis algorithms will be introduced in Section 4.3.

#### 4.1 Refining Specifications for Primary Methods

The algorithm for refinement (CA\_Gen\_Solve) is given in Fig 8. As illustrated in Section 2.2, the analysis proceeds in two steps for a primary method with shape information given in specification, namely (1) forward analysis (at lines 1-2) and (2) pure constraint abstraction generation and solving (at lines 3-10).

<pre> <b>Algorithm</b> CA.Gen.Solve(<math>\mathcal{T}, mn, e, \Phi_{pr}, \Phi_{po}, u^*, v^*</math>) 1  <math>\Delta := \text{Symb.Exec}(\mathcal{T}, mn, e, \Phi_{pr})</math> 2  <b>if</b> <math>\Delta = \text{fail}</math> <b>then return fail end if</b> 3  Normalise <math>\Delta</math> to DNF, and denote as <math>\bigvee_{i=1}^m \Delta_i</math> 4  <math>w^* := \{u^*, v^*, v'^*\} \cup \text{pureV}(\{u^*, v^*, v'^*\}, \Phi_{pr} \vee \Phi_{po})</math> 5  <math>\Delta_p := \text{Pure\_CA\_Gen}(\Phi_{po}, Q(w^*) := \bigvee_{i=1}^m \Delta_i)</math> 6  <b>if</b> <math>\Delta_p = \text{fail}</math> <b>then return fail end if</b> 7  <math>\pi := \text{Pure\_CA\_Solve}(P(w^*) := \Delta_p)</math> 8  <math>R := t\ mn\ ((t\ u)^*; (t\ v)^*)\ \text{requires}</math>        <math>\text{ex\_quan}(\Phi_{pr}, \pi)\ \text{ensures}\ \text{ex\_quan}(\Phi_{po}, \pi)</math> 9  <b>if</b> <math>\text{Verify}(\mathcal{T}, mn, R)</math> <b>then return</b> <math>\mathcal{T} \cup \{R\} \setminus</math>        <math>\{t\ mn\ ((t\ u)^*; (t\ v)^*)\ \text{requires}\ \Phi_{pr}\ \text{ensures}\ \Phi_{po}\}</math> 10 <b>else return fail end if</b> <b>end Algorithm</b> </pre>	<pre> <b>Algorithm</b> Symb.Exec (<math>\mathcal{T}, mn, e, \Phi_{pr}</math>) 11 <math>errLbbs := \emptyset</math> 12 <b>do</b> 13  <math>(\Delta, l) := \llbracket e \rrbracket_{\mathcal{T}}^{mn}(\Phi_{pr}, 0)</math> 14  <b>if</b> <math>l &gt; 0 \wedge l \notin errLbbs</math> <b>then</b> 15    <math>\Phi_{pr} := \text{ex\_quan}(\Phi_{pr}, \Delta)</math>; 16    <math>errLbbs := errLbbs \cup \{l\}</math> 17  <b>else if</b> <math>l &gt; 0 \wedge l \in errLbbs</math>        <b>then return fail</b> 18  <b>end if</b> 19 <b>while</b> <math>l &gt; 0</math> 20 <b>return</b> <math>\Delta</math> <b>end Algorithm</b> </pre>
--	---

**Fig. 8.** Refining method specifications.



The forward analysis is captured as algorithm `Symb_Exec` to the right of Fig 8. Starting from a given pre-shape  $\Phi_{pr}$ , it analyses the method body  $e$  to compute the post-state in constraint abstraction form. The symbolic execution rules are given in the appendix. They are similar to symbolic rules used in [35, 39], except for a novel mechanism to derive pure precondition, which we refer to as *pure abduction*.

This pure abduction mechanism is invoked whenever symbolic execution fails to prove memory safety based on the current prestate. For example, if we have  $ll(x, n)$  as the current state and we require  $x \mapsto \text{node}(\_, p)$  to update the value of  $p$ , then it will fail as  $ll(x, n)$  does not necessarily guarantee  $x \mapsto \text{node}(\_, p)$ . In this case we conduct the pure abduction as

$$ll(x, n) \wedge [n \geq 1] \triangleright x \mapsto \text{node}(\_, p) * \text{true}$$

to compute the missing pure information (in the squared bracket) such that the LHS (including the newly gained pure part) entails the RHS. The variable *errLbIs* (initialised at line 11) is to record the program locations in which previous pure abductions occurred. Whenever the symbolic execution fails, it returns a state  $\Delta$  that contains the pure abduction result and the location  $l$  where failure was detected, as shown in line 13. If the current abduction location  $l$  is not recorded in *errLbIs*, it indicates that this is a new failure. The abduction result is added to the precondition of the current method to obtain a stronger  $\Phi_{pr}$ , before the algorithm enters the symbolic execution loop with variable *errLbIs* updated to add in the new failure location  $l$ . This loop is repeated until symbolic execution succeeds with no memory error, or a previous failure point was re-encountered. The latter indicates either a program bug or a specification error. For example, for a method `void foo (...)` {`node w := new node(0, null); goo(w); ...`} invoking a method `goo(x)` with precondition  $ll(x, n) \wedge n \geq 2$ , our analysis will perform an abduction to get  $n \geq 2$  since it is not implied by the current state. However, as  $n$  is for the shape of local variable  $w$ , it will be quantified away when propagating  $n \geq 2$  back, ending up with `true` being added to `foo`'s precondition. In the next round of symbolic execution, our analysis will have the same abduction at the same point.

Back to the main algorithm `CA_Gen_Solve`, the analysis next builds a heap-based constraint abstraction mechanism, named  $Q(w^*)$ , for the post-state in steps 3-4. This constraint abstraction is possibly recursive. (Definition  $\dagger$  in page 5 is an example of this heap-based abstraction.) We then make use of another algorithm in Fig 9, named `Pure_CA_Gen`, to extract a pure constraint abstraction, named  $P(w^*)$ , without any heap property. (Definition  $\ddagger$  in page 6 is an instance of this pure abstraction.) This algorithm tries to derive a branch  $P_i$  for each branch  $\Delta_i$  of  $Q$ . For every  $\Delta_i$  it proceeds in two steps. In the first step (lines 22-24), it replaces the recursive occurrence of  $Q$  in  $\Delta_i$  with  $\sigma * P(w^*)$ . In the second step (lines 25-26) it tries to derive  $P_i$  via the entailment. If the entailment fails, then pure abduction is used to discover any missing pure constraint  $\sigma'_i$  for  $\rho \Delta_i$  to allow the entailment to succeed. In this case,  $\sigma'_i$  is incorporated into  $\sigma_i$  (and eventually  $P_i$ ). Once this is done, we use some existing fixpoint analysis (e.g. [38]) inside

Pure\_CA\_Solve to derive non-recursive constraint  $\pi$ , as a simplification of  $P(w^*)$ . This result is then incorporated into the pre/post specifications in line 8, before we perform a post verification in line 9 using the HIP verifier [35], to ensure the strengthened precondition is strong enough for memory safety.

Two auxiliary functions used in the algorithm are described here. The function  $\text{pureV}(V, \Delta)$  retrieves from  $\Delta$  the shapes referred to by all pointer variables from  $V$ , and returns the set of logical variables used to record numerical (size and bag) properties in these shapes, e.g.  $\text{pureV}(\{x\}, \text{ll}(x, n))$  returns  $\{n\}$ . This function is used in the algorithm to ensure that all free variables in  $\Phi_{pr}$  and  $\Phi_{po}$  are added into the parameter list of the constraint abstraction  $Q$ . The function  $\text{ex\_quan}(\Delta, \pi)$  is to strengthen the state  $\Delta$  with the abduction result  $\pi$ :  $\text{ex\_quan}(\Delta, \pi) =_{df} \Delta \wedge \exists(\text{fv}(\pi) \setminus \text{fv}(\Delta)) \cdot \pi$ . It is used to incorporate the discovered missing pure constraints into the original specification. For example,  $\text{ex\_quan}(\text{ll}(x, n), 0 < m \wedge m \leq n)$  returns  $\text{ll}(x, n) \wedge 0 < n$ .

<p><b>Algorithm</b> Pure_CA_Gen(<math>\sigma, Q(w^*) ::= \bigvee_{i=1}^m \Delta_i</math>)</p> <p>21 <b>for</b> <math>i = 1</math> <b>to</b> <math>m</math></p> <p>22   Denote all appearances of <math>Q(w^*)</math> in <math>\Delta_i</math> as <math>Q_j(w_j^*), j = 1, \dots, p</math></p> <p>23   Denote substitutions <math>\rho_j = [[w_j^*/w^*]\sigma * P(w_j^*)]/Q_j(w_j^*)]</math></p> <p>24   Let substitution <math>\rho := \rho_1 \circ \rho_2 \circ \dots \circ \rho_p</math> as applying all substitutions defined above in sequence</p> <p>25   <b>if</b> <math>(\rho\Delta_i \vdash \sigma * \sigma_i</math> <b>or</b> <math>\rho\Delta_i \wedge [\sigma'_i] \triangleright \sigma * \sigma_i)</math> <b>and</b> <math>\text{ispure}(\sigma_i)</math> <b>then</b> <math>P_i := \sigma_i</math></p> <p>26   <b>else return fail end if</b></p> <p>27 <b>end for</b></p> <p>28 <b>return</b> <math>\bigvee_{i=1}^m P_i</math></p> <p><b>end Algorithm</b></p>
---

**Fig. 9.** Pure constraint abstraction generation algorithm.

## 4.2 Pure abduction mechanism

We use the SLEEK prover [35] to check  $\Delta_1$  entails  $\Delta_2$ . If the entailment holds it can also calculate the frame  $\Delta_3$  such that  $\Delta_1 \vdash \Delta_2 * \Delta_3$ . However, when an entailment checking fails, we assume that the user has supplied necessary shape information in the specifications for primary methods, and accordingly use our pure abduction mechanism (Fig. 10) to discover missing pure constraints so that the entailment relationship becomes valid.

Our pure abduction deals with three different cases. The first rule (**R1**) applies when the LHS ( $\sigma$ ) does not entail the RHS ( $\sigma_1$ ) but the RHS entails the LHS with some pure formula ( $\sigma'$ ) as the frame; e.g. in  $\text{ll}(x, n) \not\vdash x \mapsto \text{node}(\_, \text{null})$ , the RHS can entail the LHS with pure frame  $n=1$ . The abduction then checks to

$$\boxed{
\begin{array}{c}
\frac{\sigma \not\models \sigma_1 * \mathbf{true} \quad \sigma_1 \vdash \sigma * \sigma' \quad \mathit{ispure}(\sigma') \quad \sigma \wedge \sigma' \vdash \sigma_1 * \sigma_2}{\sigma \wedge [\sigma'] \triangleright \sigma_1 * \sigma_2} \quad (\mathbf{R1}) \\
\\
\frac{\sigma \not\models \sigma_1 * \mathbf{true} \quad \sigma_1 \not\models \sigma * \mathbf{true} \quad \sigma_0 \in \mathbf{unroll}(\sigma) \quad \mathbf{data\_no}(\sigma_0) \leq \mathbf{data\_no}(\sigma_1) \quad (\sigma_0 \vdash \sigma_1 * \sigma' \text{ or } \sigma_0 \wedge [\sigma'_0] \triangleright \sigma_1 * \sigma') \quad \mathit{ispure}(\sigma') \quad \sigma \wedge \sigma' \vdash \sigma_1 * \sigma_2}{\sigma \wedge [\sigma'] \triangleright \sigma_1 * \sigma_2} \quad (\mathbf{R2}) \\
\\
\frac{\sigma \not\models \sigma_1 * \mathbf{true} \quad \sigma_1 \not\models \sigma * \mathbf{true} \quad \sigma_1 \wedge [\sigma'_1] \triangleright \sigma * \sigma' \quad \mathit{ispure}(\sigma') \quad \sigma \wedge \sigma' \vdash \sigma_1 * \sigma_2}{\sigma \wedge [\sigma'] \triangleright \sigma_1 * \sigma_2} \quad (\mathbf{R3})
\end{array}
}$$

Fig. 10. Pure abduction rules.

ensure  $\mathbf{ll}(\mathbf{x}, \mathbf{n}) \wedge \mathbf{n}=1 \vdash \mathbf{x} \mapsto \mathbf{node}(\_, \mathbf{null}) * \sigma_2$  for some  $\sigma_2$ , and returns the result  $\mathbf{n}=1$ . Note the check  $\mathit{ispure}(\sigma')$  ensures that  $\sigma'$  contains no heap information.

In the second rule (R2), neither side entails the other but the LHS term could be unfolded. An example is  $\sigma = \mathbf{sllB}(\mathbf{x}, \mathbf{S})$ ,  $\sigma_1 = \mathbf{x} \mapsto \mathbf{node}(\mathbf{u}, \mathbf{p}) * \mathbf{p} \mapsto \mathbf{node}(\mathbf{v}, \mathbf{null})$ . As the shape predicates in the antecedent are formed by disjunctions according to their definitions (like the  $\mathbf{sllB}$ ), certain branches of  $\sigma$  may entail  $\sigma_1$ . As the rule suggests, to accomplish abduction  $\sigma \wedge [\sigma'] \triangleright \sigma_1 * \sigma_2$ , we first unfold  $\sigma$  and try entailment or further abduction with the results ( $\sigma_0$ ) against  $\sigma_1$ . If it succeeds with a pure frame  $\sigma'$ , then we confirm the abduction by checking  $\sigma \wedge \sigma' \vdash \sigma_1 * \sigma_2$ . For the example above, the abduction returns  $|\mathbf{S}|=2$  ( $\sigma'$ ) and discovers the non-trivial frame  $\mathbf{S}=\{\mathbf{u}, \mathbf{v}\} \wedge \mathbf{u} \leq \mathbf{v}$  ( $\sigma_2$ ). Note that function  $\mathbf{data\_no}$  returns the number of data nodes in a state, e.g. it returns one for  $\mathbf{x} \mapsto \mathbf{node}(\mathbf{v}, \mathbf{p}) * \mathbf{ll}(\mathbf{p}, \mathbf{m})$ . (This syntactic check is important for the termination of the abduction.) The  $\mathbf{unroll}$  operation unfolds all shape predicates once in  $\sigma$ , normalises the result to a disjunctive form ( $\bigvee_{i=1}^u \sigma^i$ ), and returns the result as a set of formulae ( $\{\sigma^1, \dots, \sigma^u\}$ ). An instance is that it expands  $\mathbf{x} \mapsto \mathbf{node}(\mathbf{v}, \mathbf{p}) * \mathbf{ll}(\mathbf{p}, \mathbf{m})$  to be  $\{\mathbf{x} \mapsto \mathbf{node}(\mathbf{v}, \mathbf{p}) \wedge \mathbf{p} = \mathbf{null} \wedge \mathbf{m}=0, \exists \mathbf{u}, \mathbf{q}, \mathbf{k} \cdot \mathbf{x} \mapsto \mathbf{node}(\mathbf{v}, \mathbf{p}) * \mathbf{p} \mapsto \mathbf{node}(\mathbf{u}, \mathbf{q}) * \mathbf{ll}(\mathbf{q}, \mathbf{q}) \wedge \mathbf{m}=\mathbf{k}+1\}$ .

In the third rule (R3), neither side entails the other and the LHS term cannot be unfolded. e.g.,  $\sigma = \mathbf{x} \mapsto \mathbf{node}(\mathbf{u}, \mathbf{p}) * \mathbf{p} \mapsto \mathbf{node}(\mathbf{v}, \mathbf{null})$ ,  $\sigma_1 = \exists \mathbf{S} \cdot \mathbf{sllB}(\mathbf{x}, \mathbf{S})$ . In this case, our rule reverses the two sides of the entailment and applies the second rule to uncover the pure constraints  $\sigma'_1$  and  $\sigma'$ . It checks that adding  $\sigma'$  to the LHS ( $\sigma$ ) does entail the RHS ( $\sigma_1$ ) before it returns  $\sigma'$ . For the example above, the abduction returns  $\mathbf{u} \leq \mathbf{v}$  which is essential for the two nodes to form a sorted list (RHS).

### 4.3 Inferring Specifications for Auxiliary Methods and Loops

For auxiliary methods, we conduct a pre-analysis (Fig 11) to synthesise the pre- and post-shapes before we conduct the refinement analysis from Fig 8. Loops are dealt with by analysing their tail-recursive versions in the same way. as we do not expect the user to provide specification annotations, This approach alleviates the

need for users to provide specification annotations for both loops and auxiliary methods.

<b>Algorithm SynPre</b> $(S, f, u^*, v^*, \sigma, x^*, y^*)$ 1 $C := \text{ShpCand}(S, u^*, v^*)$ 2 <b>for</b> $\sigma_C \in C$ <b>do</b> 3 <b>if</b> $\sigma \not\vdash [x^*/u^*, y^*/v^*]\sigma_C$ 4 <b>then</b> $C := C \setminus \{\sigma_C\}$ <b>end if</b> 5 <b>end for</b> 6 <b>return</b> $C$ <b>end Algorithm</b>	<b>Algorithm SynPost</b> $(\mathcal{T}, S, f, e, \Phi_{pr}, u^*, v^*)$ 7 $C := \text{ShpCand}(S, u^*, v^*)$ 8 $\mathcal{T}' := \mathcal{T} \cup \{f(u^*, v^*) \text{ requires } \Phi_{pr} \text{ ensures false } \{e\}\}$ 9 $\Delta := \text{Symb\_Exec}(\mathcal{T}', f, \text{syn\_unroll}(f, e), \Phi_{pr})$ 10 <b>for</b> $\sigma_C \in C$ <b>do</b> 11 <b>if</b> $\Delta \wedge [\sigma] \not\vdash \sigma_C$ <b>then</b> $C := C \setminus \{\sigma_C\}$ <b>end if</b> 12 <b>end for</b> 13 <b>return</b> $\text{pair\_spec\_list}(\Phi_{pr}, C)$ <b>end Algorithm</b>
--	--

**Fig. 11.** Shape synthesis algorithms.

The pre-shape synthesis algorithm **SynPre** (Fig 11 left) takes in as input the set of shape predicates ( $S$ ), the auxiliary method name ( $f$ ), its formal parameters ( $u^*, v^*$ ), the current symbolic state in which  $f$  is called ( $\sigma$ ), and the corresponding actual parameters ( $x^*, y^*$ ) of the invocation. The algorithm first obtains possible shape candidates from the parameters  $u^*, v^*$  with **ShpCand** (line 1), then picks up a sound abstraction for the method's pre-shape with entailment, and filter out the ones which fail (line 4). Finally the pre-shape abstraction is returned. While we use an enumeration strategy here, the number of possible shape candidates per type is small as it is strictly limited by what the user provides in the primary methods, and further filtered and prioritised by our system.

To synthesise post-shapes (**SynPost**, Fig 11 right), we also assign  $C$  as possible shape candidates (line 7). We unroll  $f$ 's body  $e$  once (i.e. replace recursive calls to  $f$  in  $e$  with a substituted  $e$ ) and symbolically execute it (line 9), assuming  $f$  has a specification *requires*  $\Phi_{pr}$  *ensures* **false** (line 8). The postcondition **false** is used to ensure that the execution only considers the effect of the program branches with no recursive calls (to  $f$  itself). We then use  $\Delta$  to find out appropriate abstraction of post-shape (line 11), which is paired with  $\Phi_{pr}$  and returned as result. The function  $\text{pair\_spec\_list}(\Phi_{pr}, C)$  forms an ordered list of pre-/post-shape pairs, each of which has  $\Phi_{pr}$  as pre-shape and a  $\Phi_{po}$  in  $C$  as post-shape.

We illustrate our procedure to generate and confirm candidate shape abstractions (**ShpCand**) with an example. If we have two parameters  $x$  and  $y$  with type **node**, and the user has defined two shape predicates **l1B** and **s11B** with **node**, then the list of all possible shape candidates for the two variables ( $C$ ) will be  $[\text{s11B}(x, S) * \text{s11B}(y, T), \text{l1B}(x, S) * \text{s11B}(y, T), \text{s11B}(x, S) * \text{l1B}(y, T), \text{l1B}(x, S) * \text{l1B}(y, T), \text{s11B}(x, S), \text{s11B}(y, S), \text{l1B}(x, S), \text{l1B}(y, S), \text{emp}]$ . Elements of this list will be checked against appropriate abstract states (line 4 in Fig 11 left and line 11 in Fig 11 right) where most elements should be reduced be-

cause they are not sound abstractions. For example, in the previous list, only  $\text{llB}(\mathbf{x}, \mathbf{S}) * \text{llB}(\mathbf{y}, \mathbf{T})$  remains in the list and participates in further verification.

The initial experimental results confirm that our shape synthesis keeps only highly relevant abstractions. For the while loop in Section 2.2, we filtered out 24 (of 26) abstractions. Generally, in case that there are several abstractions as candidate specifications, we employ some other mechanisms to reduce them further. Firstly, we prioritise post-shapes with same (or stronger) predicates as in precondition since it is more likely that the output will have the same or similar shape predicates as the input, e.g.  $\mathbf{x}$  is expected to remain as  $\text{sllB}$  (or stronger) if it points to  $\text{sllB}$  as input. Secondly, we employ a lazy scheme when refining the synthesised pre/post-shapes (to complete specifications). We retrieve (and remove) the pre/post-shape pair from the head of the list, (1) use the refinement algorithm (Fig. 8) to obtain a specification for the auxiliary method, and (2) continue the analysis for the primary method. If the analysis for the primary method succeeds, we will ignore all other synthesised pre/post-shapes from the list. If either (1) or (2) fails, we will try the next one from the list. Note that our synthesis of shape specification could only cater to one predicate per parameter/result. In cases where more complex shape specifications are needed, we allow users to specify them directly for the respective auxiliary method. These mechanisms help to keep attempts over candidate specifications at a minimum level.

#### 4.4 Soundness

The underlying operational semantics of our language was given in Nguyen et al. [35]. Its concrete program state consists of stack  $s$  and heap  $h$ , as described in Section 3. The paper [35] also defined the relation  $s, h \models \Delta$  and the transitive relation  $\langle s, h, e \rangle \hookrightarrow^* \langle s', h', \nu \rangle$ . Based on that we define the soundness of our analysis as follows:

**Definition 1 (Soundness).** *For a method definition  $t \text{ mn } ((t \ u)^*; (t \ v)^*) \{e\}$ , if our analysis refines its specification as  $t \text{ mn } ((t \ u)^*; (t \ v)^*)$  requires  $\Phi_{pr}$  ensures  $\Phi_{po} \{e\}$ , then for all  $s, h \models \Phi_{pr}$ , if  $\langle s, h, e \rangle \hookrightarrow^* \langle s', h', - \rangle$ , then we have  $s', h' \models \Phi_{po}$ .*

The soundness of our analysis is ensured by the soundness of the following: the entailment prover, the abstract semantics (w.r.t. the concrete semantics), the pure constraint abstraction generation, and the fixpoint calculation. Among the above, the soundness of the entailment prover and pure fixpoint calculation are already confirmed [35, 36, 38], and hence we will concentrate on the soundness of abstract semantics and pure constraint abstraction derivation.

**Lemma 1 (Sound abstract semantics).** *If  $\llbracket e \rrbracket_{\mathcal{T}}^f(\Delta, 0) = (\Delta_1, 0)$ , then for all  $s, h$ , if  $s, h \models \Delta$  and  $\langle s, h, e \rangle \hookrightarrow \langle s_1, h_1, e_1 \rangle$ , then there always exists  $\Delta_0$  such that*

$$s_1, h_1 \models \Delta_0 \quad \text{and} \quad \llbracket e_1 \rrbracket_{\mathcal{T}}^f(\Delta_0, 0) = (\Delta_1, 0)$$

The proof is by induction and is in Appendix B.

**Lemma 2 (Sound pure constraint abstraction).** *Given a method with pre/post shape templates  $\text{pre}$  and  $\text{post}$ , if our analysis successfully computes a constraint abstraction  $Q$  in the first step, and derives a pure constraint  $P$  in the second step, then we have  $Q \vdash \text{post} \wedge P$ .*

The proof is also in Appendix B. Then based on the discussion above we have:

**Theorem 1 (Soundness).** *Our analysis is sound with respect to the underlying operational semantics.*

We have one more note about the post verification conducted in line 12 of our main algorithm in Figure 8. Such verification is used to confirm that the strengthened precondition can guarantee memory safety. This added precaution is because sometimes our refinement of precondition might not be sufficient for the program to execute without inappropriate memory access, which could be attributed to the fact that users have not provided a good enough predicate to describe the obligation of memory safety. For example, if our analysis is only supplied with a list predicate which does not contain its length information, then we can never obtain the prerequisite “the input length should be at least  $n$ ”, even if the memory safety requires that. Hence our post verification will rule out this case. However it does not affect the soundness of our analysis: the memory violation will incur **false** as an abstract state which implies any postcondition we may discover. The only reason we need it is as aforesaid — to leave only meaningful refined specifications (which has safety guarantee for the program) in our results.

## 5 Experiments and Evaluation

We have implemented a prototype system for evaluation. Our experimental results were achieved with an Intel Core 2 CPU 2.66GHz with 8Gb RAM. The four columns in Fig. 12 (in the last page) describe, resp., the analysed programs, the analysis time in seconds, and the primary methods’ (given and inferred) preconditions and postconditions. All formulae with a grey background are inferred by our analysis. We have tested 26 programs with 17 shape predicates. For some programs, we have verified them with different pre/post shape templates. Programs with star \* have different versions for various data structures.

The results highlight the refinement of both pre- and postconditions based on user-provided shape specifications, even for complicated data structures such as AVL and red-black trees. Firstly, our approach can compute non-trivial pure constraints for postcondition, e.g. for **create** we obtain the value range in the created list, for **delete** we know the content of the result list is subsumed by that of the input list, for list-sorting algorithms we confirm the content of the output is the same as that of the input, and for tree-processing programs (**insert**, **delete** and **avl\_ins**), we obtain that the height difference between the input and output trees is at most one. Meanwhile, we can calculate non-trivial requirements in precondition for memory safety or functional correctness. As an example, the

`travrs` method, taking in a list with length `m` and an integer `n`, traverses towards the tail of the list for `n` steps. the analysis discovers  $m \geq n$  in the precondition to ensure memory safety. Another example is the `append` method concatenating two sorted lists into one. To ensure that the result list is sorted, the analysis figures out that the minimum value in the second list must be no less than the maximum value in the first list.

A second highlight is our flexibility by supporting multiple predicates. Our analysis tries to refine different specifications for the same program at various correctness levels (with different predicates), e.g. `sort_insert` and `append`. For `rand_insert`, which inserts a node into a random place (after the head) of a list, we confirm that the list’s length is increased by one, but cannot verify the list is kept sorted if it was before the insertion, as the result indicates.

Another highlight is that we can reduce user annotations by synthesising specifications for auxiliary methods, given raw specifications of primary methods. For example, we have analysed a number of list-sorting algorithms with at least one auxiliary method each. We list two auxiliary methods (`merge` for `merge_sort` and `flatten` for `tree_sort`) and their discovered specifications. Note that these sorting algorithms have the same specification for their primary methods (line  $\star$ ). As another example, `avl_ins` also has some auxiliary (recursive) methods such as calculation of tree’s height, which are automatically analysed as well.

We have also tried our approach over part of the FreeRTOS kernel [1]. For its list processing programs `list.h` and `list.c` (472 lines with intensive manipulation over composite sorted doubly-linked lists) it took 2.85 seconds for our prototype to refine all the specifications given for the main functions, which further confirms the viability of our approach.

## 6 Related Work and Conclusion

### 6.1 Related works

In recent years, dramatic advances have been made in automated verification of pointer safety for heap-manipulating programs. We highlight some of them here. The local shape analysis by Distefano et al. [15] was able to infer automatically loop invariants for list-processing programs, which formed the early-version SpaceInvader tool. Gotsman et al. [17] proposed an interprocedural shape analysis for the SLayer tool. Berdine et al. [2] extended the local shape analysis [15] to handle higher-order list predicate so that more complicated real-world data structures can be analysed. Yang et al. [45] proposed a novel abstraction operation which significantly improves the scalability of the analysis. Recently, more large industrial code can be verified by the SpaceInvader tool using the compositional analysis with bi-abductive inference [5, 14].

Several shape analyses also tried to make good use of size information. In the development of the THOR tool, Magill et al. [30] proposed an adaptive shape analysis where additional numerical analysis can be used to help gain better precision. Its abstraction mechanism is also employed in C-to-gate hardware

synthesis [11]. Very recently, Magill et al. [32] formulated a novel instrumentation process which inserts numerical instructions into programs, based on their shape analysis and user-provided predicates. Instrumented programs can then be used to generate pure numerical programs for further analysis. Different from their work, we take *both* shape and numerical information into consideration when performing the abstraction, and derive the numerical abstraction from the shape constraint abstraction. Our approach can be more precise as we have more information for the abstraction. Furthermore, we can directly handle data structures with stronger invariants, like sortedness and height-balanced, which have not been addressed in THOR, to the best of our knowledge. Gulwani et al. [18] combine a set domain with its cardinality domain in a general framework. Compared with these, our approach can handle data structures with stronger invariants like sortedness, height-balanced and bag-related invariants, which have not been addressed in the previous works. Another piece of work, by Chang et al. [6] and Chang and Rival [7], employs inductive checkers and checker segments to express shape and numerical information. Our previous loop invariant synthesis [39] also infers strong loop invariants with both shape and numerical information but is limited to while loop analysis. Compared with their works, ours addresses specification refinement with pure properties (including numerical and bag ones) in both pre- and postconditions by processing shape and pure information in two phases with the help of pure abduction. Our previous loop invariant synthesis [39] also infers strong loop invariants with a one-phase heavyweight abstract interpretation. Compared with this work, it is limited to loop analysis, whereas this work tackles not only loops but also methods; meanwhile this work is more lightweight as it solves the constraint abstraction in two phases where the second phase (pure constraint abstraction solving) utilises existing provers and is hence more modular and efficient.

There are also many other approaches to expressing heap-based domains than separation logic. Hackett and Rugina [22] can deal with AVL-trees but is customised to handle only tree-like structures with height property. The shape analysis framework TVLA [41] is based on three-valued logic. It is capable of handling complicated data structures and properties, such as sortedness. LRP [46] is fully decidable over multiple linked data structures and has a finite model property. Guo et al. [19] reported a global shape analysis that discover inductive structural shape invariants from the code. Kuncak et al. [27] developed a role system to express and track referencing relationships among objects, where an object's role (type) depends on, and changes according to, the mutation of its referencing. Compared with these works, separation logic based approach benefits from the frame rule and hence supports local reasoning. Meanwhile, our approach heads towards full functional correctness including bag-related properties, which previous ones do not generally handle.

There are also numerous works on automated assertion discovery, for example those based on abstract interpretation [12]. Compared with our work, they mainly focus on finding numerical program properties, and hence our work is complementary to them in the light that we also discover heap/shape informa-



tion. Meanwhile, we can utilise such works as our pure solver, for example the disjunction inference [38].

On the verification side, Smallfoot [3] is the first verification system based on separation logic. The HIP/SLEEK verification system [34, 35] supports user-defined shape predicates over the combined shape and numerical domain. The SLEEK tool has played a very important role in our analysis. The PALE system [33] transforms constraints in the pointer assertion logic (PAL) into monadic second-order logic (MSO) and discharge them with MONA. Hob [44] is a modular program verification tool for shape properties. Based on Hob, Jahob [25, 44] takes Java as its target language and allows more general specification language. Havoc [8] is another verification tool for C language about heap-allocated data structures, using a novel reachability predicate. There is another recent work on refining specifications via counterexample-guided abstraction refinement [42] which is goal-driven and incrementally improves for given safety requirements. Among these works, our verification is distinguished because we free users from writing whole specifications by requiring only partial specifications, and omit user-supplied annotations for less important loops and auxiliary methods.

## 6.2 Conclusion

We have reported a new approach to program verification that accepts partial specifications of methods, and refines them by discovering missing constraints for numerical and bag properties, aiming at full functional correctness for pointer-based data structures. We further augment our approach by requiring only partial specification for primary methods. Specifications for loops and auxiliary methods can then be systematically discovered. We have built a prototype system and the initial experimental results are encouraging.

**Acknowledgement.** This work was supported in part by the EPSRC projects [EP/E021948/1, EP/G042322/1] and MoE ARF grant R-252-000-411-112.

## References

1. R. Barry. FreeRTOS — a free RTOS for small embedded real time systems. 2006.
2. J. Berdine, C. Calcagno, B. Cook, D. Distefano, P. O’Hearn, T. Wies, and H. Yang. Shape analysis for composite data structures. In *19th CAV*, July 2007.
3. J. Berdine, C. Calcagno, and P. O’Hearn. Smallfoot: Modular automatic assertion checking with separation logic. In *FMCO*, 2005.
4. J. Berdine, B. Cook, D. Distefano, and P. O’Hearn. Automatic termination proofs for programs with shape-shifting heaps. In *18th CAV*, 2006.
5. C. Calcagno, D. Distefano, P. O’Hearn, and H. Yang. Compositional shape analysis by means of bi-abduction. In *36th POPL*, January 2009.
6. B. Chang, X. Rival, and G. Nacula. Shape analysis with structural invariant checkers. In *SAS*, 2007.
7. B. Chang and X. Rival. Relational inductive shape analysis. In *POPL*, 2008.
8. S. Chatterjee, S. Lahiri, S. Qadeer, and Z. Rakamaric. A reachability predicate for analyzing low-level software. In *TACAS*, 2007.
9. W.-N. Chin, C. David, H. H. Nguyen, and S. Qin. Automated verification of shape, size and bag properties. In *12th ICECCS*, 2007.

10. W.-N. Chin, C. David, H. H. Nguyen, and S. Qin. Enhancing modular oo verification with separation logic. In *POPL*, January 2008.
11. B. Cook, A. Gupta, S. Magill, A. Rybalchenko, J. Simsa, and V. Vafeiadis. Finding heap-bounds for hardware synthesis. In *FMCAD*, November 2009.
12. P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *POPL*, 1977.
13. P. Cousot and R. Cousot. On abstraction in software verification. In *CAV*, 2002.
14. D. Distefano. Attacking large industrial code with bi-abductive inference. In *FMICS*, 2009.
15. D. Distefano, P. O'Hearn, and H. Yang. A local shape analysis based on separation logic. In *TACAS*, 2006.
16. R. Giacobazzi. Abductive analysis of modular logic programs. In *ILPS*, 1994.
17. A. Gotsman, J. Berdine, and B. Cook. Interprocedural shape analysis with separated heap abstractions. In *SAS*, 2006.
18. S. Gulwani, T. Lev-Ami, and M. Sagiv. A Combination Framework for Tracking Partition Sizes. In *POPL*, 2009.
19. B. Guo, N. Vachharajani, and D. August. Shape analysis with inductive recursion synthesis. In *PLDI*, 2007.
20. A. Gupta, R. Majumdar, and A. Rybalchenko. From tests to proofs. In *TACAS*, 2009.
21. J. Gustavsson and J. Svenningsson. Constraint abstractions. In *Programs as Data Objects II*, Denmark, May 2001.
22. B. Hackett and R. Rugina. Region-based shape analysis with tracked locations. In *POPL*, 2005.
23. J. Henriksen, J. Jensen, M. Jørgensen, N. Klarlund, R. Paige, T. Rauhe, and A. Sandholm. Mona: Monadic second-order logic in practice. In *TACAS*, 1995.
24. S. Ishtiaq and P. O'Hearn. BI as an assertion language for mutable data structures. In *POPL*, 2001.
25. V. Kuncak. *Modular Data Structure Verification*. PhD thesis, EECS Department, Massachusetts Institute of Technology, February 2007.
26. S. Nejati. *Refinement Relation on Partial Specifications*. MSc thesis, CS Department, University of Toronto, 2003.
27. V. Kuncak, P. Lam, and M. Rinard. Role analysis. In *POPL*, 2002.
28. V. Kuncak, P. Lam, K. Zee, and M. Rinard. Modular pluggable analyses for data structure consistency. *IEEE Transactions on Software Engineering*, 32(12), 2006.
29. S. Magill, A. Nanevski, E. Clarke, and P. Lee. Inferring Invariants in Separation Logic for Imperative List-processing Programs. In *the SPACE Workshop*, 2006.
30. S. Magill, J. Berdine, E. Clarke, and B. Cook. Arithmetic strengthening for shape analysis. In *SAS*, 2007.
31. S. Magill, M. Tsai, P. Lee, and Y. Tsay. Thor: A tool for reasoning about shape and arithmetic. In *CAV*, 2008.
32. S. Magill, M. Tsai, P. Lee, and Y. Tsay. Automatic numeric abstractions for heap-manipulating programs. In *POPL*, 2010. To appear.
33. A. Møller and M. Schwartzbach. The pointer assertion logic engine. *ACM SIGPLAN Notices*, 36(5):221–231, 2001.
34. H. H. Nguyen and W.-N. Chin. Enhancing program verification with lemmas. In *20th CAV*, 2008.
35. H. H. Nguyen, C. David, S. Qin, and W.-N. Chin. Automated verification of shape and size properties via separation logic. In *8th VMCAI*, 2007.
36. T. Nipkow, L. C. Paulson, and M. Wenzel. Isabelle/HOL — a proof assistant for higher-order logic, volume 2283 of LNCS. Springer, 2002.
37. M. Parkinson and G. Bierman. Separation logic, abstraction and inheritance. In *POPL*, 2008.
38. C. Popeea and W.-N. Chin. Inferring disjunctive postconditions. In *Proceedings of 11th Asian Computing Science Conference*, 2006.

39. S. Qin, G. He, C. Luo, and W.-N. Chin. Loop invariant synthesis in a combined domain. In *ICFEM*, 2010. To appear.
40. J. Reynolds. Separation logic: a logic for shared mutable data structures. In *17th LICS*, 2002.
41. M. Sagiv, T. Reps, and R. Wilhelm. Parametric shape analysis via 3-valued logic. *ACM Transactions on Programming Languages and Systems*, 24(3):217–298, 2002.
42. M. Taghdiri. *Automating Modular Verification by Refining Specifications*. PhD thesis, EECS Department, Massachusetts Institute of Technology, February 2008.
43. J. Warford. *Computer systems*. Jones & Bartlett Publishers, 2009.
44. T. Wies, V. Kuncak, K. Zee, A. Podelski, and M. Rinard. On verifying complex properties using symbolic shape analysis. *The Computing Research Repository*, abs/cs/0609104, 2006.
45. H. Yang, O. Lee, J. Berdine, C. Calcagno, B. Cook, D. Distefano, and P. O’Hearn. Scalable shape analysis for systems code. In *20th CAV*, 2008.
46. G. Yorsh, A. Rabinovich, M. Sagiv, A. Meyer, and A. Bouajjani. A logic of reachable patterns in linked data-structures. In *FOSSACS*, 2006.

## A Symbolic Execution Rules

This section defines the symbolic execution rules used in the first step of the constraint abstraction generation. If the program contains recursive calls to itself, the postcondition will be in a recursive (open) form.

The type of our symbolic execution is defined as

$$\llbracket e \rrbracket =_{df} \text{AllSpec} \rightarrow \text{Names} \rightarrow (\mathcal{P}_{\text{SH}} \times \text{Int}) \rightarrow (\mathcal{P}_{\text{SH}} \times \text{Int})$$

where **AllSpec** contains all the specifications of all methods (extracted from the program *Prog*) and **Names** contains all the method names. The integer (label) in both input and output is used to record a program location where abduction is needed. If the integer remains zero after the symbolic execution of  $e$ , then the output state denotes the post-state of  $e$ . However, a positive number indicates that an abduction must have occurred and the resulting state (the abduction result) will be propagated back to the the method’s precondition by our analysis, so that the next round of symbolic execution should succeed in the same location.

The foundation of the symbolic execution is the basic transition functions from a conjunctive abstract state to a conjunctive or disjunctive abstract state below:

$$\begin{aligned} \text{unfold}(x) &=_{df} \text{SH} \rightarrow \mathcal{P}_{\text{SH}[x]} && \text{Unfolding} \\ \text{exec}(d[x]) &=_{df} (\text{AllSpec} \times \text{Names}) \rightarrow (\text{SH}[x] \times \text{Int}) \rightarrow (\text{SH} \times \text{Int}) && \text{Heap-sensitive exec.} \\ \text{exec}(d) &=_{df} (\text{AllSpec} \times \text{Names}) \rightarrow (\text{SH} \times \text{Int}) \rightarrow (\text{SH} \times \text{Int}) && \text{Heap-insensitive exec.} \end{aligned}$$

where  $\text{SH}[x]$  denotes the set of conjunctive abstract states in which each element has  $x$  exposed as the head of a data node ( $c(x, v^*)$ ), and  $\mathcal{P}_{\text{SH}[x]}$  contains all the (disjunctive) abstract states, each of which is composed by such conjunctive states. Here  $\text{unfold}(x)$  unfolds the symbolic heap so that the cell referred to by  $x$  is exposed for access by heap sensitive commands  $d[x]$  via the second transition function  $\text{exec}(d[x])$ . The third function defined for other (heap insensitive) commands  $d$  does not require such exposure of  $x$ .

For the unfolding operation  $\text{unfold}(x)$ , there are two possible scenarios. If  $x$  refers to a data node in the current state  $\sigma$ , no unfolding is required and the  $\text{exec}$  operation can proceed directly. However, if  $x$  refers to a (user-defined) shape predicate, then  $\text{unfold}(x)$  will unfold the current state  $\sigma$  according to the definition of the predicate in order to expose the data node referred to by  $x$ :

$$\frac{\text{isdatat}(c) \quad \sigma \vdash c(x, v^*) * \sigma'}{\text{unfold}(x)\sigma \rightsquigarrow \sigma} \quad \frac{\text{isspred}(c) \quad \sigma \vdash c(x, u^*) * \sigma' \quad c(\text{root}, v^*) \equiv \Phi}{\text{unfold}(x)\sigma \rightsquigarrow \sigma' * [x/\text{root}, u^*/v^*]\Phi}$$

The test  $\text{isdatat}(c)$  returns **true** only if  $c$  is a data node and  $\text{isspred}(c)$  returns **true** only if  $c$  is a shape predicate.

The symbolic execution of heap-sensitive commands  $d[x]$  (i.e.  $x.f$ ,  $x.f := w$ , or  $\text{free}(x)$ ) assumes that the unfolding  $\text{unfold}(x)$  has been done prior to the execution. The first three rules below are for normal symbolic execution where the current state is sufficiently strong for safe execution. The last two rules handle the cases where the symbolic execution fails and abductive reasoning can be used to discover missing pure information.

$$\begin{array}{c} \frac{\text{isdatat}(c) \quad \sigma \vdash c(x, v_1, \dots, v_n) * \sigma'}{\text{exec}(x.f_i)(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow (\sigma' * c(x, v_1, \dots, v_n) \wedge \text{res} = v_i, 0)} \\ \frac{\text{isdatat}(c) \quad \sigma \vdash c(x, v_1, \dots, v_n) * \sigma'}{\text{exec}(x.f_i := w)(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow (\sigma' * c(x, v_1, \dots, v_{i-1}, w, v_{i+1}, \dots, v_n), 0)} \\ \frac{\text{isdatat}(c) \quad \sigma \vdash c(x, u^*) * \sigma'}{\text{exec}(\text{free}(x))(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow (\sigma', 0)} \\ \frac{\text{isdatat}(c) \quad \sigma \not\vdash c(x, u^*) * \text{true} \quad \sigma * [\sigma'] \triangleright c(x, u^*) * \text{true}}{\text{exec}(d[x])(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow (\sigma', \text{lbl}(d[x]))} \\ \frac{\text{isdatat}(c) \quad \sigma \not\vdash c(x, u^*) * \text{true} \quad \sigma * [\sigma'] \not\triangleright c(x, u^*) * \text{true}}{\text{exec}(d[x])(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow (\text{false}, \text{lbl}(d[x]))} \end{array}$$

Note that the second to last rule uses an abductive reasoning (via SLEEK) to discover the missing numerical information  $\sigma'$ . Here we use a mapping  $\text{lbl}(-)$  to map any instruction in the program being analysed to a unique positive integer label (namely the aforementioned program location). The rule changes the second element of the result to  $\text{lbl}(d[x])$  which will be used by the analysis to record the instruction causing an abduction, quits the current execution, propagates the discovered information back to the precondition of the current method, and restarts the symbolic execution with the strengthened precondition. The last rule covers the scenario in which the abduction fails. Then the execution cannot continue and returns  $(\text{false}, \text{lbl}(d[x]))$ .

The symbolic execution rules for heap-insensitive commands are as follows:

$$\begin{array}{c}
\text{exec}(k)(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow (\sigma \wedge \text{res}=k, 0) \quad \text{exec}(v)(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow (\sigma \wedge \text{res}=v, 0) \\
\\
\frac{\text{isdatat}(c)}{\text{exec}(\text{new } c(v^*))(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow (\sigma * c(\text{res}, v^*), 0)} \\
\\
\frac{(x^*, y^*) = \text{vars}(w, e) \quad (g, \mathcal{T}_1) = \text{Analysis}(\mathcal{T}, \text{while}(w)\{e\}, \sigma, x^*, y^*) \quad \mathcal{T}' = \mathcal{T} \cup \mathcal{T}_1}{\text{exec}(\text{while}(w)\{e\})(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow \text{exec}(f(x^*; y^*))(\mathcal{T}', g)(\sigma, 0)} \\
\\
\frac{t \text{ mn } ((t_i \ u_i)_{i=1}^m; (t_i \ v_i)_{i=1}^n) \in \mathcal{T} \quad (mn, \mathcal{T}_1) = \text{Analysis}(\mathcal{T}, mn, \sigma, x^*, y^*) \quad \mathcal{T}' = \mathcal{T} \cup \mathcal{T}_1}{\text{exec}(mn(x_1..x_m; y_1..y_n))(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow \text{exec}(mn(x_1..x_m; y_1..y_n))(\mathcal{T}', mn)(\sigma, 0)} \\
\\
\frac{t \text{ mn } ((t_i \ u_i)_{i=1}^m; (t_i \ v_i)_{i=1}^n) \text{ requires } \Phi_{pr} \text{ ensures } \Phi_{po} \in \mathcal{T} \quad mn \neq f}{\begin{array}{c} \rho = [x'_i/u'_i]_{i=1}^m \circ [y'_i/v'_i]_{i=1}^n \quad \sigma \vdash \rho \Phi_{pr} * \sigma' \quad \rho_o = [y_i/v_i]_{i=1}^n \circ \rho \\ \rho_l = [r_i/y'_i]_{i=1}^n \quad \rho_{ol} = [r_i/y_i]_{i=1}^n \quad \text{fresh logical } r_i \end{array}}{\text{exec}(mn(x_1..x_m; y_1..y_n))(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow ((\rho_l \sigma') * (\rho_{ol} \circ \rho_o \Phi_{po}), 0)} \\
\\
\frac{t \text{ mn } ((t_i \ u_i)_{i=1}^m; (t_i \ v_i)_{i=1}^n) \text{ requires } \Phi_{pr} \text{ ensures } \Phi_{po} \in \mathcal{T} \quad mn = f}{\begin{array}{c} \rho = [x'_i/u'_i]_{i=1}^m \circ [y'_i/v'_i]_{i=1}^n \quad \sigma \vdash \rho \Phi_{pr} * \sigma' \quad \rho_o = [y_i/v_i]_{i=1}^n \circ \rho \\ \rho_l = [r_i/y'_i]_{i=1}^n \quad \rho_{ol} = [r_i/y_i]_{i=1}^n \quad \text{fresh logical } r_i \end{array}}{\text{exec}(mn(x_1..x_m; y_1..y_n))(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow ((\rho_l \sigma') * (\rho_{ol} \circ \rho_o (\Phi_{po} \wedge \text{P}(u^*, v^*)))), 0)} \\
\\
\frac{t \text{ mn} \dots \in \mathcal{T} \quad \rho = [x'_i/u'_i]_{i=1}^m \circ [y'_i/v'_i]_{i=1}^n \quad \sigma \not\vdash \rho \Phi_{pr} * \text{true} \quad \sigma \wedge [\sigma'] \triangleright \rho \Phi_{pr} * \text{true}}{\text{exec}(mn(x_1..x_m; y_1..y_n))(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow (\sigma', \text{lbl}(mn(\dots)))} \\
\\
\frac{t \text{ mn} \dots \in \mathcal{T} \quad \rho = [x'_i/u'_i]_{i=1}^m \circ [y'_i/v'_i]_{i=1}^n \quad \sigma \not\vdash \rho \Phi_{pr} * \text{true} \quad \sigma \wedge [\sigma'] \not\triangleright \rho \Phi_{pr} * \text{true}}{\text{exec}(mn(x_1..x_m; y_1..y_n))(\mathcal{T}, f)(\sigma, 0) \rightsquigarrow (\text{false}, \text{lbl}(mn(\dots)))}
\end{array}$$

Note that the first three rules deal with constant ( $k$ ), variable ( $v$ ) and data node creation ( $\text{new } c(v^*)$ ), respectively, while the remaining rules handle method invocation. The fourth and fifth rules are used for the invocation of a while loop or an auxiliary method which has not been analysed, where we employ the analysis algorithm recursively to achieve its postcondition to enable application of the next rule. The sixth rule is used for the invocation of another method  $mn$  which has already been analysed ( $mn \neq f$ ), and the call site meets the precondition of  $mn$ , as checked by the entailment  $\sigma \vdash \rho \Phi_{pr} * \sigma'$ . In this case, the execution succeeds and moves on. The seventh rule is for a recursive call to the current method ( $mn = f$ ), similar as above except that a constraint abstraction is in place as postcondition. The last two rules are for the cases where the call site cannot establish the precondition of the callee method and where abductive reasoning is employed. In both cases, the execution discontinues. The eighth rule returns the abduction result  $\sigma'$ , which is a pure formula and will be propagated back by the analysis to strengthen the caller method's precondition. The last rule captures the scenario in which the abduction fails. Note that the operator  $\circ$  is used to compose two substitutions: the substitution  $\rho_2 \circ \rho_1$  works by first applying  $\rho_1$  and then  $\rho_2$ .

To keep presentation simple, we assume there are no mutual recursions in the programs to analyse; therefore each method to be analysed should only call itself recursively. This assumption does not lose generality, as we can always transform mutual recursion into single recursion [43] to have only one constraint abstraction  $Q$  in our analysis for one method.

The following rule for all commands signifies that when starting from a configuration in which the second element is positive (i.e. a faulty state), the execution will not change the state. This rule is used to skip all remaining instructions when abductive reasoning is used as a new round of symbolic execution with strengthened precondition should be started instead:

$$\frac{l > 0}{\text{exec}(-)(\mathcal{T}, f)(\sigma, l) \rightsquigarrow (\sigma, l)}$$

We can now lift  $\text{unfold}$ 's domain to  $\mathcal{P}_{\text{SH}}$  using the following operation  $\text{unfold}^\dagger$ :

$$\text{unfold}^\dagger(x) \bigvee \sigma_i =_{df} \bigvee (\text{unfold}(x) \sigma_i)$$

and similarly for  $\text{exec}$ :

$$\text{exec}^\dagger(d)(\mathcal{T}, f)(\bigvee \sigma_i, l) =_{df} (\bigvee \sigma'_i, \max\{l_i\}) \text{ where } (\sigma'_i, l_i) \in \text{exec}(d)(\mathcal{T}, f)(\sigma_i, l)$$

The symbolic execution rules for program constructors  $e$  can now be defined using the lifted transition functions above. Firstly, no change will be made if starting from a faulty state, as the first rule shows. In all other cases, the symbolic execution transforms one abstract state to another w.r.t. the program instruction:

$$\begin{aligned} \llbracket - \rrbracket_{\mathcal{T}}^f(\Delta, l) &=_{df} (\Delta, l), \text{ where } l > 0 \\ \llbracket d[x] \rrbracket_{\mathcal{T}}^f(\Delta, 0) &=_{df} \text{exec}^\dagger(d[x])(\mathcal{T}, f)(\text{unfold}^\dagger(x) \Delta, 0) \\ \llbracket d \rrbracket_{\mathcal{T}}^f(\Delta, 0) &=_{df} \text{exec}^\dagger(d)(\mathcal{T}, f)(\Delta, 0) \\ \llbracket e_1; e_2 \rrbracket_{\mathcal{T}}^f(\Delta, 0) &=_{df} \llbracket e_2 \rrbracket_{\mathcal{T}}^f \circ \llbracket e_1 \rrbracket_{\mathcal{T}}^f(\Delta, 0) \\ \llbracket v := e \rrbracket_{\mathcal{T}}^f(\Delta, 0) &=_{df} [v_1/v', r_1/\text{res}](\llbracket e \rrbracket_{\mathcal{T}}^f(\Delta, 0)) \wedge v' = r_1, \text{ fresh } v_1, r_1 \\ \frac{(\Delta'_1, l_1) = \llbracket e_1 \rrbracket_{\mathcal{T}}^f(v \wedge \Delta, 0) \quad (\Delta'_2, l_2) = \llbracket e_2 \rrbracket_{\mathcal{T}}^f(\neg v \wedge \Delta, 0)}{\llbracket \text{if } (v) \ e_1 \ \text{else } e_2 \rrbracket_{\mathcal{T}}^f(\Delta, 0) =_{df} (\Delta'_1 \vee \Delta'_2, \max\{l_1, l_2\})} \end{aligned}$$

## B Soundness

This section proves the lemmas and theorem for soundness in Section 4.

**Lemma 1 (Sound abstract semantics).** *If  $\llbracket e \rrbracket_{\mathcal{T}}^f(\Delta, 0) = (\Delta_1, 0)$ , then for all  $s, h$ , if  $s, h \models \Delta$  and  $\langle s, h, e \rangle \hookrightarrow \langle s_1, h_1, e_1 \rangle$ , then there always exists  $\Delta_0$  such that*

$$s_1, h_1 \models \Delta_0 \quad \text{and} \quad \llbracket e_1 \rrbracket_{\mathcal{T}}^f(\Delta_0, 0) = (\Delta_1, 0)$$

**Proof** The proof is done by structural induction over program constructors:

- Case **null** |  $k$  |  $v$  |  $v.f$ . Straightforward.
- Case  $v := e$ . There are two cases according to the operational semantics:
  - $e$  is not a value. From operational semantics, there is  $e_1$  s.t.  $\langle s, h, e \rangle \hookrightarrow \langle s_1, h_1, e_1 \rangle$ , and  $\langle s, h, v := e \rangle \hookrightarrow \langle s_1, h_1, v := e_1 \rangle$ . From abstract semantics for assignment, if  $\llbracket e \rrbracket_{\mathcal{T}}^f(\Delta, 0) = (\Delta_2, 0)$ , and  $\Delta_1 = [v_1/v', r_1/\text{res}](\Delta_2) \wedge v' = r_1$ . By induction hypothesis, there exists  $\Delta_0, s_1, h_1 \models \Delta_0$  and  $\llbracket e_1 \rrbracket_{\mathcal{T}}^f(\Delta_0, 0) = (\Delta_2, 0)$ . It concludes from the assignment rule that  $\llbracket v := e_1 \rrbracket_{\mathcal{T}}^f(\Delta_0, 0) = (\Delta_1, 0)$ .
  - $e$  is a value. Trivial.
- Case  $v_1.f := v_2$ . Take  $S_0 = S$ . It concludes immediately from the **exec** rule for field update and the underlying operational semantics.
- Case **new**  $c(v^*)$ . From abstract semantics for **new**, we have  $\llbracket \text{new } c(v^*) \rrbracket_{\mathcal{T}}^f(\Delta, 0) = (\Delta_1, 0)$ , where  $\Delta_1 = \Delta * c(\text{res}, v'_1, \dots, v'_n)$ . Let  $\Delta_0 = \Delta_1$ . From the operational semantics, we have  $\langle s, h, \text{new } c(v^*) \rangle \hookrightarrow \langle s, h + [\iota \mapsto r], \iota \rangle$ , where  $\iota \notin \text{dom}(h)$ . From  $s, h \models \Delta$ , we have  $s, h + [\iota \mapsto r] \models \Delta_0$ . Moreover,  $\llbracket \iota \rrbracket_{\mathcal{T}}^f(\Delta_0, 0) = (\Delta_1, 0)$ .
- Case  $e_1; e_2$ . We consider the case where  $e_1$  is not a value (otherwise it is straightforward). From the operational semantics, we have  $\langle s, h, e_1 \rangle \hookrightarrow \langle s_1, h_1, e_3 \rangle$ . From the abstract semantics rule for sequence, we have  $\vdash \{\Delta\} e_1 \{ \Delta_2 \}$ . By induction hypothesis, there exists  $\Delta_0$  s.t.  $s_1, h_1 \models \text{Post}(\Delta_0)$ , and  $\vdash \{\Delta_0\} e_3 \{ \Delta_2 \}$ . By the sequential rule we have  $\llbracket e_3; e_2 \rrbracket_{\mathcal{T}}^f(\Delta_0, 0) = \Delta_1$ .
- Case **if**  $(v)$   $e_1$  **else**  $e_2$ . There are two possibilities in the operational semantics:
  - $s(v) = \text{true}$ . We have  $\langle s, h, \text{if } (v) e_1 \text{ else } e_2 \rangle \hookrightarrow \langle s, h, e_1 \rangle$ . Let  $\Delta_0 = (\Delta \wedge v')$ . It is obvious that  $s, h \models \Delta_0$ . From the if-conditional rule of abstract semantics, we have:

$$\begin{aligned} \llbracket e_1 \rrbracket_{\mathcal{T}}^f(\Delta_0, 0) &= (\Delta_2, 0) \\ \llbracket e_2 \rrbracket_{\mathcal{T}}^f(\Delta \wedge \neg v', 0) &= (\Delta_3, 0) \end{aligned}$$

And we also have (due to sound weakening of postcondition)

$$\llbracket e_1 \rrbracket_{\mathcal{T}}^f(\Delta_0, 0) = (\Delta_2 \vee \Delta_3, 0)$$

That is,  $\llbracket e_1 \rrbracket_{\mathcal{T}}^f(\Delta_0, 0) = (\Delta_1, 0)$ .

- $s(v) = \text{false}$ . Analogous.
- Case  $mn(v_{1..n})$ . For the method invocation rule, we know  $\Delta \vdash [v'_j/v_j]_{j=1}^n \Phi_{pr}^i * \Delta^i$ , for  $i = 1, \dots, p$ . Take  $\Delta_0 = \bigvee_{i=1}^p [v'_j/v_j]_{j=1}^n \Phi_{pr}^i * \Delta^i$ . From the operational semantics and the above heap entailment, we have  $s_1, h_1 \models \Delta_0$ . Then the method invocation rule implies  $\forall i \in 1..p \cdot \llbracket e_1 \rrbracket_{\mathcal{T}}^f([v'_j/v_j]_{j=1}^n \Phi_{pr}^i * \Delta^i, 0) = (\Delta^i * \Phi_{po}^i, 0)$ . Therefore we have  $\llbracket e_1 \rrbracket_{\mathcal{T}}^f(\Delta_0, 0) = (\Delta_1, 0)$  which concludes.
- Case **while**  $(v)$   $\{e\}$ . It can be converted to tail-recursive method call with all parameters passed by reference, and thus follows the above case.  $\square$

**Lemma 2 (Sound pure constraint abstraction).** *Given a method with pre/post shape templates `pre` and `post`, if our analysis successfully computes a constraint abstraction  $Q$  in the first step, and derives a pure constraint  $P$  in the second step, then we have  $Q \vdash \text{post} \wedge P$ .*

**Proof** This proof follows directly our procedure to compute the pure constraint abstraction from the shape one. Denote the shape constraint abstraction as

$$Q ::= \bigvee_i Q_i$$

and the provided post-shape `post`, we use

$$Q_i \vdash \text{post} \wedge P_i$$

to derive each  $P_i$ , and construct

$$P ::= \bigvee_i P_i$$

Therefore, our result  $Q \vdash \text{post} \wedge P$  can be obtained from the fact that  $\bigvee_i Q_i \vdash \text{post} \wedge (\bigvee_i P_i)$ .  $\square$

**Theorem 2 (Soundness).** *Our analysis is sound with respect to the underlying operational semantics.*

**Proof** For a method  $t \text{ mn } ((t \ u)^*; (t \ v)^*) \text{ requires } \Phi_{pr} \text{ ensures } \Phi_{po}\{e\}$ , suppose we have the resulted pure constraint  $\pi$ . Then, for all  $s, h \models \Phi_{pr} \wedge \text{ex\_quan}(\Phi_{pr}, \pi)$  and  $\llbracket e \rrbracket_{\mathcal{T}}^{mn}(\Phi_{pr} \wedge \text{ex\_quan}(\Phi_{pr}, \pi), 0) = (\Delta, 0)$ , we know  $s', h' \models \Delta$  (Lemma 1). Thus let  $\Delta = \bigvee Q_i$ , we have  $s', h' \models \text{lfix } \lambda Q. \bigvee Q_i$  where  $\text{lfix}$  is the least fixed-point operator. Then due to Lemma 2 and the soundness of entailment, we know  $s', h' \models \Phi_{po} \wedge \text{lfix } \lambda P. \bigvee P_i$ . Because  $\pi = \text{lfix } \lambda P. \bigvee P_i$ , we can claim  $s', h' \models \Phi_{po} \wedge \pi$ . Finally from the definition of  $\text{ex\_quan}(\Phi_{po}, \pi)$  (as its result is a weakening of  $\pi$  by quantifying out its local variables) we know  $s', h' \models \Phi_{po} \wedge \text{ex\_quan}(\Phi_{po}, \pi)$ .  $\square$



Prog.	Time	Pre	Post
List processing programs			
create*	0.379	$\text{emp} \wedge n \geq 0$	$\text{llB}(\text{res}, S) \wedge n =  S  \wedge \forall v \in S. 1 \leq v \leq n$
	1.752	$\text{emp} \wedge n \geq 0$	$\text{dllB}(\text{res}, \text{rp}, S) \wedge n =  S  \wedge \forall v \in S. 1 \leq v \leq n$
	0.954	$\text{emp} \wedge n \geq 0$	$\text{sllB2}(\text{res}, S) \wedge n =  S  \wedge \forall v \in S. 1 \leq v \leq n$
sort_* insert	0.591	$\text{ll}(x, n) \wedge n \geq 1$	$\text{ll}(x, m) \wedge m = n + 1$
	0.789	$\text{dll}(x, p, n) \wedge n \geq 1$	$\text{dll}(x, q, m) \wedge n \geq 1 \wedge m = n + 1 \wedge p = q$
	0.504	$\text{sll}(x, n, \text{xs}, \text{x1}) \wedge v \geq \text{xs}$	$\text{sll}(x, m, \text{mn}, \text{mx}) \wedge \text{xs} = \text{mn} \wedge \text{mx} = \max(\text{x1}, v) \wedge m = n + 1$
tail_ insert	0.566	$\text{ll}(x, n) \wedge n \geq 1$	$\text{ll}(x, m) \wedge m = n + 1$
	0.628	$\text{sll}(x, n, \text{xs}, \text{x1}) \wedge v \geq \text{x1}$	$\text{sll}(x, m, \text{mn}, \text{mx}) \wedge v = \text{mx} \wedge \text{mn} = \text{xs} \wedge m = n + 1$
rand_* insert	0.522	$\text{ll}(x, n) \wedge n \geq 1$	$\text{ll}(x, m) \wedge m = n + 1$
	0.830	$\text{dll}(x, p, n) \wedge n \geq 1$	$\text{dll}(x, q, m) \wedge m = n + 1 \wedge p = q$
	—	$\text{sll}(x, n, \text{xs}, \text{x1}) \wedge (\text{fail})$	$\text{sll}(x, m, \text{mn}, \text{mx}) \wedge (\text{fail})$
delete	0.630	$\text{llB}(x, S) \wedge  S  \geq 2$	$\text{llB}(x, T) \wedge \exists a. S = T \sqcup \{a\}$
	1.024	$\text{sllB}(x, S) \wedge  S  \geq 2$	$\text{sllB}(x, T) \wedge \exists a. S = T \sqcup \{a\}$
delete	1.252	$\text{dllB}(x, p, S) \wedge  S  \geq 2$	$\text{dllB}(x, q, T) \wedge \exists a. S = T \sqcup \{a\} \wedge p = q$
travrs	0.296	$\text{ll}(x, m) \wedge n \geq 0 \wedge m \geq n$	$\text{ls}(x, p, k) * \text{ll}(\text{res}, r) \wedge p = \text{res} \wedge k = n \wedge m = n + r$
	2.205	$\text{sllB}(x, S) \wedge n \geq 0 \wedge  S  \geq n$	$\text{slsB}(x, p, T) * \text{sllB}(\text{res}, S_2) \wedge p = \text{res} \wedge  T  = n$ $\wedge S = T \sqcup S_2 \wedge \forall u \in T, v \in S_2. u \leq v$
append*	0.512	$\text{ll}(x, \text{xn}) * \text{ll}(y, \text{yn}) \wedge \text{xn} \geq 1$	$\text{ll}(x, m) \wedge m = \text{xn} + \text{yn}$
	0.660	$\text{dll}(x, \text{xp}, \text{xn}) * \text{dll}(y, \text{yp}, \text{yn}) \wedge \text{xn} \geq 1$	$\text{dll}(x, q, m) \wedge m = \text{xn} + \text{yn} \wedge q = \text{xp}$
	0.948	$\text{sll}(x, \text{xn}, \text{xs}, \text{x1}) \wedge \text{x1} \leq \text{ys} * \text{sll}(y, \text{yn}, \text{ys}, \text{y1})$	$\text{sll}(x, m, \text{rs}, \text{r1}) \wedge \text{y1} = \text{r1} \wedge m \geq 1 + \text{yn} \wedge m = \text{xn} + \text{yn}$
Sorting (main)		$\text{llB}(x, S) \wedge  S  \geq 1$	$\text{sllB}(\text{res}, T) \wedge T = S \quad (*)$
merge	4.107	$\text{sllB}(x, S_x) * \text{sllB}(y, S_y)$	$\text{sllB}(\text{res}, T) \wedge T = S_x \sqcup S_y$
flatten	2.693	$\text{bstB}(x, S)$	$\text{sllB}(\text{res}, T) \wedge T = S$
insert	0.824	$\text{sllB}(r, S) * x \mapsto \text{node}(v, -)$	$\text{sllB}(\text{res}, T) \wedge T = S \sqcup \{v\}$
quick	2.132	$\text{lbd}(x, S)$	$\text{lbd}(x, S_1) * \text{lbd}(\text{res}, S_2) \wedge S = S_1 \sqcup S_2 \wedge \forall u \in S_1 \forall v \in S_2. u \leq p \leq v$
Binary tree, binary search tree, AVL tree and red-black tree processing programs			
travrs	0.532	$\text{bt}(x, S, h)$	$\text{bt}(x, T, k) \wedge S = T \wedge h = k$
count	0.709	$\text{bt}(x, S, h)$	$\text{bt}(x, T, k) \wedge \text{res} =  S  \wedge S = T \wedge h = k$
height	0.913	$\text{bt}(x, S, h)$	$\text{bt}(x, T, k) \wedge \text{res} = h = k \wedge S = T$
insert	1.276	$\text{bt}(x, S, h) \wedge  S  \geq 1 \wedge h \geq 1$	$\text{bt}(x, T, k) \wedge T = S \sqcup \{v\} \wedge h \leq k \leq h + 1$
delete	0.970	$\text{bt}(x, S, h) \wedge  S  \geq 2 \wedge h \geq 2$	$\text{bt}(x, T, k) \wedge \exists a. S = T \sqcup \{a\} \wedge h - 1 \leq k \leq h$
search	1.583	$\text{bst}(x, \text{sm}, \text{lg})$	$\text{bst}(x, \text{mn}, \text{mx}) \wedge \text{sm} = \text{mn} \wedge \text{lg} = \text{mx} \wedge 0 \leq \text{res} \leq 1$
bst_ insert	1.720	$\text{bst}(x, \text{sm}, \text{lg})$	$\text{bst}(x, \text{mn}, \text{mx}) \wedge (v < \text{sm} \wedge v = \text{mn} \wedge \text{lg} = \text{mx} \vee \text{lg} < v \wedge v = \text{mx} \wedge \text{sm} = \text{mn} \vee \text{sm} = \text{mn} \wedge \text{lg} = \text{mx})$
avl_ins	11.12	$\text{avl}(x, S, h)$	$\text{avl}(\text{res}, T, k) \wedge T = S \sqcup \{v\} \wedge h \leq k \leq h + 1$
rbt_ins	8.76	$\text{rbt}(x, S)$	$\text{rbt}(\text{res}, T) \wedge T = S \sqcup \{v\}$
sdl2nbt	5.826	$\text{sdlB}(x, p, q, S) \wedge  S  \geq 1 \wedge p = \text{null} \wedge q = \text{tail}$	$\text{nbt}(\text{res}, T) \wedge T = S$

Fig. 12. Experimental Results.