

Automated Specification Discovery via User-Defined Predicates

Guanhua He¹, Shengchao Qin^{1*}, Wei-Ngan Chin² and Florin Craciun³

¹Teeside University

²National University of Singapore

³Babes-Bolyai University

Abstract. Automated discovery of specifications for heap-manipulating programs is a challenging task due to the complexity of aliasing and mutability of data structures. This task is further complicated by an expressive domain that combines shape, numerical and bag information. In this paper, we propose a compositional analysis framework in the presence of user-defined predicates, which would derive the summary for each method in the expressive abstract domain, independently from its callers. We propose a novel abstraction method with a bi-abduction technique in the combined domain to discover pre-/post-conditions that could not be automatically inferred before. The analysis does not only prove the memory safety properties, but also finds relationships between pure and shape domains towards full functional correctness of programs. A prototype of the framework has been implemented and initial experiments have shown that our approach can discover interesting properties for non-trivial programs.

1 Introduction

In automated program analysis, certain kinds of program properties have been well explored over the last decades, such as numerical properties in linear abstraction domain, and shape properties for list-manipulating programs in separation domain. However, previous works have not yet automatically analysed program properties involving complex mixed domains, particularly for programs with sophisticated data structures and strong invariants involving both structural and pure (numerical and content) information. For example, it is still non-trivial to discover program properties, such as a list becoming sorted during the execution of a program, a binary search tree remaining balanced before and after the execution of a procedure, or the elements of a list remain unchanged after reversing the list. This difficulty is not only due to sharing and mutability of data structures under manipulation, but is also due to closely intertwined program properties, such as structural numerical information (length and height), symbolic contents of data structures (bag of values), and relational numerical information (sortedness and balancedness).

In addition to classical shape analyses (e.g. [4, 14, 24]), separation logic [22] has been applied to analyse shape properties in recent years [5, 8, 26]. These works can automatically infer method specifications in the shape domain. Some other works such as [17, 18] also incorporate simple numerical information into the shape domain to allow automated synthesis of properties like data structure size information.

However, these previous analyses mainly deal with predesignated data structure properties with fixed numerical templates, such as pointer safety for lists and list length

* Corresponding author: shengchao.qin@gmail.com

information. To overcome this limitation, we propose in this paper a compositional program analysis in a combined abstract domain with *shape*, *numerical* and *bag* information. Our analysis not only handles both functional correctness and memory safety together, but can also discover relationships between shape and pure (numerical and bag) domains. Unlike traditional approaches [18] which usually analyse the shape first before turning to pure properties, our approach analyses programs over both domains at the same time. This is very necessary as verifying functional correctness for certain programs may require us to consider both shape and pure information at the same time. Without pure information, a shape analysis may not be able to find useful program specifications (an example is the `merge` procedure discussed in [5]). Our approach can handle this kind of programs smoothly, and we will illustrate our method using the `merge` example in Section 2.

Our analysis is compositional. It analyses a program fragment without any given contextual information, and it analyses each method in a modular way independent of its callers. To generate the summary (pre-/post-conditions) for each method, our analysis adopts a new bi-abduction mechanism over the combined domain, which generalises the bi-abduction technique proposed by Calcagno et al. [5] to a more expressive abstract domain. In summary, this paper makes the following contributions:

- We have designed a compositional analysis to discover *full program specifications* (in the form of pre-/post-conditions involving shape, numerical and bag properties) with user-given data structure predicates.
- For such an analysis, we have designed a *bi-abductive abstract semantics* which incorporates a generalised bi-abduction procedure to facilitate specification discovery over the combined abstract domain.
- In addition to a normal abstraction function, we have also proposed a novel *abductive abstraction* function over the combined domain. This new abstraction function allows us to find stronger method specifications that are often necessary for the successful verification for higher level of functional correctness.
- We have built a prototype system and conducted some initial experiments, which help confirm the viability and precision of our solution in inferring non-trivial program specifications.

2 The Approach

In this section we give some preliminaries and illustrate our approach via an example.

2.1 Preliminaries

Separation Logic. Separation logic [22] extends Hoare logic to support reasoning about shared mutable data structures. It provides separation conjunction ($*$) to form formulae like $p_1 * p_2$ to assert that two heaps described by p_1 and p_2 are domain-disjoint.

User-defined Predicates. In our analysis, users are allowed to define inductive predicates in separation logic to specify both separation and pure properties of recursive data structures. For example, given a data structure `data Node { int val; Node next; }`, one can define a predicate for a list with its content as

$$\begin{aligned} \text{llB}(\text{root}, n, S) &\equiv (\text{root}=\text{null} \wedge n=0 \wedge S=\emptyset) \vee \\ &(\exists v, q, n_1, S_1 \cdot \text{root} \mapsto \text{Node}(v, q) * \text{llB}(q, n_1, S_1) \wedge n_1 = n - 1 \wedge S = S_1 \sqcup \{v\}) \end{aligned}$$

The parameter `root` for the predicate `llB` is the root pointer referring to the list. The length and content of the list are denoted resp. by `n` and the bag `S`, and \sqcup indicates multi-set (bag) union. If one wants to verify a sorting algorithm, they can specify a non-empty sorted list as follows:

$$\begin{aligned} \text{sllB}(\text{root}, \text{mi}, \text{mx}, S) &\equiv (\text{root} \mapsto \text{Node}(\text{mi}, \text{null}) \wedge \text{mi} = \text{mx} \wedge S = \{\text{mi}\}) \vee \\ &(\text{root} \mapsto \text{Node}(v, q) * \text{sllB}(q, \text{m}_1, \text{mx}, S_1) \wedge v = \text{mi} \wedge v \leq \text{m}_1 \wedge \text{m}_1 \leq \text{mx} \wedge S = S_1 \sqcup \{v\}) \end{aligned}$$

where it keeps track of the minimum (`mi`) and maximum (`mx`) values in the list as well as the bag of all values (`S`). Note that we use a shortened notation that unbound variables, such as `q`, `v`, `m1` and `S1`, are implicitly existentially quantified.

Such predicates play an important role in our analysis as (i) they are used to help specify desired properties about data structures under manipulation, and (ii) they serve as a guide for our analysis to discover desired program specifications. To reduce the burden of supplying such predicates, we have defined a library of predicates covering popular data structures and variety of properties.

Entailment. In our work we make use of the separation logic prover SLEEK [7] to prove whether one formula Δ' in the combined abstract domain entails another one Δ : $\Delta' \vdash \Delta * R$. R is called the *frame* which is useful for our analysis. For instance, by entailment proof

$$\exists y \cdot x \mapsto \text{node}(vx, y) * \text{llB}(y, n, S) \vdash \text{llB}(x, m, S_1) * R$$

We can generate the frame R as $m = n + 1 \wedge S_1 = S \sqcup \{vx\}$.

Bi-Abduction. In an earlier work [5], a bi-abductive entailment is proposed for the *shape* domain: given two shape formulae G, H , the bi-abduction $G * [A] \triangleright H * [F]$ infers the *anti-frame* A and the *frame* F along the entailment proof. An example taken from [5] is

$$x \mapsto \text{null} * z \mapsto \text{null} * [\text{list}(y)] \triangleright \text{list}(x) * \text{list}(y) * [z \mapsto \text{null}]$$

where the `list(·)` predicate describes acyclic, null-terminated singly-linked lists. In the current work, we will generalise such bi-abductive reasoning to the combined domain (involving shape, user-defined predicates, numerical and bag information). A simple example of the generalised bi-abductive reasoning is

$$\exists y \cdot x \mapsto \text{node}(vx, y) * y \mapsto \text{node}(vy, \text{null}) * [A] \triangleright \text{sllB}(x, \text{mi}, \text{mx}, S) * [F]$$

where $A \equiv (vx \leq vy)$ and $F \equiv (\text{mi} = vx \wedge \text{mx} = vy \wedge S = \{vx, vy\})$.

2.2 An Illustrative Example

We illustrate our analysis approach via the `merge` method (used in the merge-sort), which has been declared as an unverifiable example in [5], since their analysis does not keep track of data values stored in the list during their shape analysis. The method (Fig. 1) merges two sorted lists into one sorted list. Automated specification discovery for `merge` is tricky due to two facts: (1) only one input list is fully traversed; (2) both input lists are required to be sorted. For (1), if we apply the shape abduction [5], we can only discover two disjoint lists (for precondition) - one ending with null and one ending with an unknown pointer, which cannot guarantee the memory safety of the method. For (2), if an analysis cannot infer that the two input lists are sorted, it will

1	Node merge(Node x, Node y)	9	Node t = x.next;
2	{	10	x.next = merge(t, y);
3	if (x == null) {	11	return x;
4	return y;	12	} else {
5	} else if (y == null) {	13	Node t = y.next;
6	return x;	14	y.next = merge(x, t);
7	} else	15	return y;
8	if (x.val <= y.val) {	16	} }

Fig. 1. Merging two sorted lists.

not be able to discover that the output list is sorted, which will not be sufficient for one to verify the functional correctness of the enclosing merge-sort method. The two input lists being unsorted also causes the unknown pointer problem mentioned above. To overcome these difficulties, we propose a compositional analysis in a combined shape and pure domain, where program properties over the combined domain are processed at the same time during the analysis. Our analysis adopts a novel bi-abduction mechanism to help discover program preconditions in the combined domain.

For the `merge` example, the shape predicate selected for our analysis is `slsB` which keeps track of the minimal (`mi`) and maximal (`mx`) values, bag of values (`S`) and tail pointer (`p`) of a sorted list segment.

$$\begin{aligned} \text{slsB}(\text{root}, \text{mi}, \text{mx}, \text{S}, \text{p}) &\equiv (\text{root} \mapsto \text{Node}(\text{mi}, \text{p}) \wedge \text{mi} = \text{mx} \wedge \text{S} = \{\text{mi}\}) \vee \\ &(\text{root} \mapsto \text{Node}(\text{mi}, \text{q}) * \text{sllB}(\text{q}, \text{m}_1, \text{mx}, \text{S}_1, \text{p}) \wedge \text{mi} \leq \text{m}_1 \wedge \text{m}_1 \leq \text{mx} \wedge \text{S} = \text{S}_1 \sqcup \{\text{mi}\}) \end{aligned}$$

Our analysis aims at finding a sound and precise specification (summary) of the method. Starting from an initial specification ($\text{Pre}_0 \equiv \text{emp}$, $\text{Post}_0 \equiv \text{false}$), our analysis iterates the method body by symbolic execution a number of times until a fixed point is reached for the pre-/post-condition pair. During the symbolic execution, we use a pair of states (infP , Curr) to keep track of the precondition that the analysis has discovered (infP) so far and the current state the execution has reached (Curr), respectively. If the current abstract state does not meet the precondition required by the current program command, we use an abductive inference mechanism (mentioned in the previous subsection) to synthesise a candidate precondition as the missing precondition.

For the `merge` example, the initial specification ($\text{Pre}_0 \equiv \text{emp}$, $\text{Post}_0 \equiv \text{false}$) allows the analysis to skip the branches with recursive calls to `merge`. The symbolic execution in the first fixpoint iteration starts from state ($\text{infP} \equiv \text{emp}$, $\text{Curr} \equiv \text{emp}$), since the analysis assumes no prior knowledge about the starting program state. To enter line 4, the condition `x==null` needs to be met by the current abstract state. We apply abduction and discover `x=null` which is then added to the precondition. Similarly, we have `y=null` from the second branch. After the first iteration, a summary is found as

$$(\text{Pre}_1 \equiv (\text{x=null} \vee \text{y=null}), \text{Post}_1 \equiv (\text{x=null} \wedge \text{res=y} \vee \text{y=null} \wedge \text{res=x})) \quad (1)$$

where `res` denotes the return value. Using this new summary for recursive calls to `merge`, symbolically executing the method body again (but with an updated starting state ($\text{infP} \equiv \text{Pre}_1$, $\text{Curr} \equiv \text{Pre}_1$)) yields the summary (Pre_2 , Post_2):

$$\begin{aligned} \text{Pre}_2 &\equiv \text{x=null} \vee \text{y=null} \vee \text{x} \mapsto \text{Node}(\text{xv}_1, \text{xp}_1) * \text{y} \mapsto \text{Node}(\text{yv}_1, \text{yp}_1) \\ &\quad \wedge (\text{xv}_1 \leq \text{yv}_1 \wedge \text{xp}_1 = \text{null} \vee \text{xv}_1 > \text{yv}_1 \wedge \text{yp}_1 = \text{null}), \\ \text{Post}_2 &\equiv \text{x=null} \wedge \text{res=y} \vee \text{y=null} \wedge \text{res=x} \vee \text{x} \mapsto \text{Node}(\text{xv}_1, \text{xp}_1) * \text{y} \mapsto \text{Node}(\text{yv}_1, \text{yp}_1) \\ &\quad \wedge (\text{xv}_1 \leq \text{yv}_1 \wedge \text{res=x} \wedge \text{xp}_1 = \text{y} \vee \text{xv}_1 > \text{yv}_1 \wedge \text{res=y} \wedge \text{yp}_1 = \text{x}) \end{aligned}$$

After the third iteration of symbolic execution, we generate a precondition as: (2)

$$\begin{aligned} & x = \text{null} \vee y = \text{null} \vee x \mapsto \text{Node}(xv_1, xp_1) * y \mapsto \text{Node}(yv_1, yp_1) \\ & \wedge (xv_1 \leq yv_1 \wedge xp_1 = \text{null} \vee xv_1 > yv_1 \wedge yp_1 = \text{null}) \end{aligned} \quad (3)$$

$$\begin{aligned} & \forall x \mapsto \text{Node}(xv_1, xp_1) * xp_1 \mapsto \text{Node}(xv_2, xp_2) * y \mapsto \text{Node}(yv_1, yp_1) \\ & \wedge (xv_1 \leq yv_1 \wedge (xv_2 \leq yv_1 \wedge xp_2 = \text{null} \vee xv_2 > yv_1 \wedge yp_1 = \text{null})) \end{aligned} \quad (4)$$

$$\begin{aligned} & \forall x \mapsto \text{Node}(xv_1, xp_1) * y \mapsto \text{Node}(yv_1, yp_1) * yp_1 \mapsto \text{Node}(yv_2, yp_2) \\ & \wedge (xv_1 > yv_1 \wedge (xv_1 \leq yv_2 \wedge xp_1 = \text{null} \vee xv_1 > yv_2 \wedge yp_2 = \text{null})) \end{aligned} \quad (5)$$

Branch (4) says that the program only touches the second node of x list (the list referred to by x) if $xv_1 \leq yv_1$. Furthermore, if $xv_2 \leq yv_1$, xp_2 should be null ; otherwise yp_1 must be null to guarantee the termination of the method and memory safety. Branch (5) states a similar condition when touching the second node of y list. The information kept in this formula is very precise, but keeping such a level of details will not allow the analysis to scale up. According to the given predicate s1sB , we could abstract the shape of the x list (and that of the y list) to be a sorted list segment. However, the formula itself does not contain sufficient information for us to carry out this abstraction, i.e. the sortedness information about the x list (and the y list) is missing. This missing information is the numerical relation between xv_1 and xv_2 in the x list (and that between yv_1 and yv_2 in the y list). In other words, we need to use abduction to discover $xv_1 \leq xv_2$ (resp. $yv_1 \leq yv_2$) during the abstraction from the shape of the x list (resp. the y list) to a sorted list segment in the branch (4) (resp. (5)), e.g. for the x list:

$$x \mapsto \text{Node}(xv_1, xp_1) * xp_1 \mapsto \text{Node}(xv_2, xp_2) * [xv_1 \leq xv_2] \triangleright \text{s1sB}(x, xv_1, xv_2, xS_1, xp_2) * R$$

The inspiration for this *abductive abstraction* comes from the definition of the predicate s1sB . We use such predicates to help infer data structure properties that are anticipated from some program code. Note that a standard abstraction would only be able to obtain an abstraction of an ordinary list segment without any sortedness information.

By applying such an *abductive abstraction* against the predicate s1sB and then joining the branches with the same shape, the precondition from two iterations becomes:

$$\begin{aligned} & x = \text{null} \vee y = \text{null} \vee \text{s1sB}(x, xmi_0, xmx_0, xS_0, xp_0) * \text{s1sB}(y, ymi_0, ymx_0, yS_0, yp_0) \\ & \wedge (xmx_0 \leq ymx_0 \wedge xp_0 = \text{null} \vee xmx_0 > ymx_0 \wedge yp_0 = \text{null}) \end{aligned}$$

Continuing the analysis, the fixed point of the program summary (Pre,Post) is reached:

$$\begin{aligned} \text{Pre} & \equiv x = \text{null} \vee y = \text{null} \vee \text{s1sB}(x, xmi_0, xmx_0, xS_0, xp_0) * \\ & \text{s1sB}(y, ymi_0, ymx_0, yS_0, yp_0) \wedge (xmx_0 \leq ymx_0 \wedge xp_0 = \text{null} \vee xmx_0 > ymx_0 \wedge yp_0 = \text{null}), \\ \text{Post} & \equiv x = \text{null} \wedge \text{res} = y \vee y = \text{null} \wedge \text{res} = x \vee \text{s1sB}(x, xmi_1, xmx_1, xS_1, xp_1) \\ & * \text{s1sB}(y, ymi_1, ymx_1, yS_1, yp_1) \wedge xS_0 \sqcup yS_0 = xS_1 \sqcup yS_1 \wedge xmi_1 = xmi_0 \wedge ymi_1 = ymi_0 \wedge \\ & (xmi_0 \leq ymi_0 \wedge \text{res} = x \wedge xp_1 = y \wedge xmx_1 \leq ymi_1 \vee xmi_0 > ymi_0 \wedge \text{res} = y \wedge yp_1 = x \wedge ymx_1 \leq xmi_1) \end{aligned}$$

The essential steps to terminate the search for suitable preconditions are abstraction and widening. Both operators are tantamount to weakening a state, and they are over-approximations and are sound for the synthesis of postconditions. However, when such steps are applied to the synthesis of preconditions, it may make the precondition too weak for the program to establish the postcondition. So after the analysis, we shall use a forward analysis process to check the discovered summary (a similar process is also carried out in [5]).

From this example, we observe that the memory safety is not only related to the shape of data structures, but may also relate to data values stored in them. For the merge example, our analysis can find that one input list is traversed to its end, i.e. until null is reached, and the other input list is partially traversed till it reaches an element

that is larger than the maximal value of the former list. As captured in the inferred precondition, the rest of the list will not be accessed by the program. Similarly, the inferred postcondition captures a fairly precise specification that represents the merged list using two list segments that either begins from x or from y , depending on which of the two input lists contains the smallest element.

3 Language and Abstract Domain

To simplify presentation, we employ a strongly-typed C-like imperative language in Fig. 2 to demonstrate our approach. A program $Prog$ written in this language consists of declarations $tdecl$, which can be data type declarations $datat$ (e.g. Node in Section 2), predicate definitions $spred$ (e.g. `llB` and `s1sB`), as well as method declarations $meth$. The definitions for $spred$ and $mspec$ are given later in Fig. 3. We assume that methods come with no specifications (i.e. no $mspec^*$ part), and our proposed analysis will discover them. Our language is expression-oriented, and thus the body of a method (e) is an expression formed by program constructors. Note that d and $d[v]$ represent respectively heap-insensitive and heap sensitive commands. k^τ is a constant of type τ . The language allows both call-by-value and call-by-reference method parameters, separated with a semicolon (;). These parameters allow each iterative loop to be directly converted to an equivalent tail-recursive method, where mutations on parameters are made visible to the caller via pass-by-reference. This technique of translating away iterative loops is standard and is helpful in further minimising our core language.

$Prog ::= tdecl^* meth^*$	$tdecl ::= datat \mid spred$
$datat ::= data\ c\ \{ field^* \}$	$field ::= t\ v \quad t ::= c \mid \tau$
$meth ::= t\ mn\ ((t\ v)^*; (t\ v)^*)\ mspec^*\ \{e\}$	$\tau ::= int \mid bool \mid void$
$e ::= d \mid d[v] \mid v := e \mid e_1; e_2 \mid t\ v; e \mid if\ (v)\ e_1\ else\ e_2$	
$d ::= null \mid k^\tau \mid v \mid new\ c(v^*) \mid mn(u^*; v^*)$	
$d[v] ::= v.f \mid v_1.f := v_2 \mid free(v)$	

Fig. 2. A Core (C-like) Imperative Language.

Our specification language (in Fig. 3) allows (user-defined) shape predicates $spred$ to specify program properties in our combined domain. Note that such predicates are constructed with disjunctive constraints Φ . We require that the predicates be well-formed [7]. The first parameter of a predicate is the pointer referring to the data structures itself. A conjunctive abstract program state σ has mainly two parts: the heap (shape) part κ in the separation domain and the pure part π in convex polyhedra domain and bag (multi-set) domain, where π consists of γ , ϕ and φ as aliasing, numerical and multi-set information, respectively. k^{int} is an integer constant. The square symbols like \square , \sqsubseteq , \sqcup and \sqcap are multi-set operators. The set of all σ formulae is denoted as SH (*symbolic heap*). During the symbolic execution, the abstract program state at each program point will be a disjunction of σ 's, denoted by Δ . Its set is defined as \mathcal{P}_{SH} . An abstract state Δ can be normalised to the Φ form [7].

Using entailment [7], we define a partial order over these abstract states:

$$\Delta \preceq \Delta' =_{df} \Delta' \vdash \Delta * R$$

$spre$	$::= p(\text{root}, v^*) \equiv \Phi$	Φ	$::= \bigvee \sigma^*$	σ	$::= \exists v^*. \kappa \wedge \pi$
$mspec$	$::= \text{requires } \Phi_{pr} \text{ ensures } \Phi_{po}$				
Δ	$::= \Phi \mid \Delta_1 \vee \Delta_2 \mid \Delta \wedge \pi \mid \Delta_1 * \Delta_2 \mid \exists v. \Delta$				
κ	$::= \text{emp} \mid v \mapsto c(v^*) \mid p(v^*) \mid \kappa_1 * \kappa_2$			π	
γ	$::= v_1 = v_2 \mid v = \text{null} \mid v_1 \neq v_2 \mid v \neq \text{null} \mid \gamma_1 \wedge \gamma_2$				
ϕ	$::= \varphi \mid b \mid a \mid \phi_1 \wedge \phi_2 \mid \phi_1 \vee \phi_2 \mid \neg \phi \mid \exists v. \phi \mid \forall v. \phi$				
b	$::= \text{true} \mid \text{false} \mid v \mid b_1 = b_2$			a	
s	$::= k^{\text{int}} \mid v \mid k^{\text{int}} \times s \mid s_1 + s_2 \mid -s \mid \max(s_1, s_2) \mid \min(s_1, s_2) \mid B $				
φ	$::= v \in B \mid B_1 = B_2 \mid B_1 \sqsubset B_2 \mid B_1 \sqsubseteq B_2 \mid \forall v \in B. \phi \mid \exists v \in B. \phi$				
B	$::= B_1 \sqcup B_2 \mid B_1 \sqcap B_2 \mid B_1 - B_2 \mid \emptyset \mid \{v\}$				

Fig. 3. The Specification Language.

where R is the (computed) residue part. And we also have an induced lattice over these states as the base of fixpoint calculation for our analysis.

The memory model of our specification formulae can be found in [7]. In our analysis, variables include both program and logical variables.

4 Generalised Bi-Abduction for the Combined Domain

We present a new bi-abduction procedure over the combined domain (which generalises the previous bi-abduction [5] over only the shape domain).

Given σ and σ_1 , the bi-abduction procedure $\sigma * [\sigma'] \triangleright \sigma_1 * \sigma_2$ (shown in Fig. 4) aims to find the anti-frame part σ' and the frame part σ_2 such that $\sigma * \sigma' \vdash \sigma_1 * \sigma_2$ where σ and σ_1 can be the current program state and the precondition of next instruction, respectively. Our abduction procedure can handle more than one predicates in the analysis, while the shape abduction [5] caters for only one specified shape predicate domain. Another advance is that we can infer numerical and bag properties together with the shape formulae as the anti-frame to improve the precision of the analysis.

$$\begin{array}{c}
\frac{\sigma \not\vdash \sigma_1 * \text{true} \quad \sigma_1 \vdash \sigma * \sigma' \quad \sigma * \sigma' \vdash \sigma_1 * \sigma_2}{\sigma * [\sigma'] \triangleright \sigma_1 * \sigma_2} \text{Residue} \\
\frac{\sigma \not\vdash \sigma_1 * \text{true} \quad \sigma_1 \not\vdash \sigma * \text{true} \quad \sigma_0 \in \text{unroll}(\sigma) \quad \text{data_no}(\sigma_0) \leq \text{data_no}(\sigma_1) \quad \sigma_0 \vdash \sigma_1 * \sigma' \text{ or } \sigma_0 * [\sigma'_0] \triangleright \sigma_1 * \sigma' \quad \sigma * \sigma' \vdash \sigma_1 * \sigma_2}{\sigma * [\sigma'] \triangleright \sigma_1 * \sigma_2} \text{Unroll} \\
\frac{\sigma \not\vdash \sigma_1 * \text{true} \quad \sigma_1 \not\vdash \sigma * \text{true} \quad \sigma_1 * [\sigma'_1] \triangleright \sigma * \sigma' \quad \sigma * \sigma' \vdash \sigma_1 * \sigma_2}{\sigma * [\sigma'] \triangleright \sigma_1 * \sigma_2} \text{Reverse} \\
\frac{\sigma \not\vdash \sigma_1 * \text{true} \quad \sigma_1 \not\vdash \sigma * \text{true} \quad \sigma * \sigma_1 \vdash \sigma_1 * \sigma_2}{\sigma * [\sigma_1] \triangleright \sigma_1 * \sigma_2} \text{Missing} \\
\frac{\sigma \not\vdash \sigma_1 * \sigma'_1 * \text{true} \quad \sigma_1 * \sigma'_1 \not\vdash \sigma * \text{true} \quad \sigma \vdash \sigma'_1 * \text{true} \quad \sigma * [\sigma'] \triangleright \sigma_1 * \sigma'_2 \quad \sigma * \sigma' \vdash \sigma_1 * \sigma'_1 * \sigma_2}{\sigma * [\sigma'] \triangleright (\sigma_1 * \sigma'_1) * \sigma_2} \text{Remove}
\end{array}$$

Fig. 4. Bi-Abduction rules.

The 1st rule **Residue** triggers when the LHS (σ) does not entail the RHS (σ_1) but the RHS entails the LHS with some formula (σ') as the residue. This rule is quite general and applies in many cases. For instance, if LHS is $\text{emp}(\sigma)$, RHS is $x \mapsto \text{Node}(xv, xp)(\sigma_1)$, the RHS can entail the LHS with residue $x \mapsto \text{Node}(xv, xp)(\sigma')$. The abduction then

checks whether σ plus the frame σ' implies $\sigma_1 * \sigma_2$ for some σ_2 (emp in this example), and returns $x \mapsto \text{Node}(xv, xp)$ as the anti-frame.

The 2nd rule `Unroll` deals with the case where neither side entails the other, e.g. for `s1sB(x, xmi, xmx, xS, null)` as LHS and $\exists p, u, v. x \mapsto \text{Node}(u, p) * p \mapsto \text{Node}(v, \text{null})$ as RHS. As the shape predicates in the antecedent σ are formed by disjunctions according to their definitions (like `s1sB`), its certain disjunctive branches may imply σ_1 . As the rule suggests, to accomplish abduction $\sigma * [\sigma'] \triangleright \sigma_1 * \sigma_2$, we first unfold σ ($\sigma_0 \in \text{unroll}(\sigma)$) and try entailment or further abduction with the results (σ_0) against σ_1 . If it succeeds with a frame σ' , then we confirm the abduction by ensuring $\sigma * \sigma' \vdash \sigma_1 * \sigma_2$. For the example above, the abduction returns $\exists u, v. xS = \{u, v\}$ as the anti-frame σ' and discovers the nontrivial frame $u = xmi \wedge v = xmx \wedge u \leq v$ as σ_2 . The function `data_no` returns the number of data nodes in a state, e.g. it returns 1 for $x \mapsto \text{Node}(v, p) * \text{11B}(p, n, T)$. This syntactic check prevents unlimited number of times of unrolling from happening when the abduction procedure invokes this rule recursively. The `unroll` unfolds all shape predicates once in σ , normalises the result to a disjunctive form ($\bigvee_{i=1}^n \sigma_i$), and returns the result as a set of formulae ($\{\sigma_1, \dots, \sigma_n\}$).

The 3rd rule `Reverse` handles the case where neither side entails the other, and the 2nd rule does not apply, e.g. $\exists p, u, v, q. x \mapsto \text{Node}(u, p) * p \mapsto \text{Node}(v, q)$ as LHS and $\exists xS. s1sB(x, xmi, xmx, xS, xp)$ as RHS. In this case the antecedent cannot be unfolded as it contains only data nodes. As the rule suggests, it reverses two sides of the entailment and applies the second rule to uncover the constraints σ'_1 and σ' . Then it checks that the LHS (σ), with σ' added, does entail the RHS (σ_1) before it returns σ' . For the example above, the anti-frame is inferred as $u \leq v$.

When an abduction procedure is conducted, the first three rules should be attempted exhaustively in the given order; if they do not succeed in finding a solution, then the rule `Missing` is invoked to add the consequence to the antecedent, provided that they are consistent. It is effective for situations like $x \mapsto \text{Node}(-, -) \not\vdash y \mapsto \text{Node}(-, -)$, where we should add $y \mapsto \text{node}(-, -)$ to the LHS directly. In our analysis, we assume that different variables refer to different nodes unless aliasing is suggested in the program code. For example, the if-statement `if (x == y){c}` suggests that `x` and `y` are aliased in code `c`. Note that when the third rule is applied, the abduction procedure in the premise, namely $\sigma_1 * [\sigma'] \triangleright \sigma * \sigma'$, is not allowed to apply the third rule again. This is to prevent an infinite number of applications of the third rule.

If the first four rules fail, the `Remove` rule then tries to find a part of consequent (σ'_1) which is entailed by the antecedent. The abduction is then applied to the remaining part of the consequent (σ_1) to discover the anti-frame (σ'). For example, the bi-abduction question $\text{11B}(x, n, S) \wedge n > 2 * [\sigma'] \triangleright x \mapsto \text{Node}(v_1, p_1) * y \mapsto \text{Node}(v_2, p_2) * \sigma_2$ needs this rule to remove $x \mapsto \text{Node}(v_1, p_1)$ from consequent before applying the `Missing` rule to find the anti-frame $\sigma' = y \mapsto \text{Node}(v_2, p_2)$.

Our earlier work [20] gives a restricted form of abduction focusing on discovering pure information with the assumption that either complete or partial shape information is available. Our bi-abduction algorithm presented here generalises it to cater for full specification discovery scenarios, whereby, we do not have the hints to guide the analysis anymore due to the absence of shape information in pre/post-conditions; but at the same time we can have more freedom as to what missing information to discover.

One observation on abduction is that there can be many solutions of the anti-frame σ' for the entailment $\sigma * \sigma' \vdash \sigma_1 * \sigma_2$ to succeed. Therefore, we define “quality” of anti-frame solutions with the partial order \preceq given in the previous section, i.e. the smaller (weaker) one is regarded as better. We prefer to find solutions that are (potentially locally) minimal with respect to \preceq and consistent. However, such solutions are generally not easy to compute and could incur excess cost (with additional disjunction in the analysis). Therefore, our abductive inference is designed more from a practical perspective to discover anti-frames that should be suitable as preconditions for programs, and the partial order \preceq sounds more like a guidance of the decision choices of our abduction implementation, rather than a guarantee to find the theoretically best solution.

5 Analysis Algorithm

Our proposed analysis algorithm is given in Fig. 5. It takes three input parameters: \mathcal{T} as the set of method specifications that are already inferred, the procedure to be analysed $t\ mn\ ((t\ x)^*; (t\ y)^*)\ \{e\}$, and a pre-set upper bound n on the number of shared logical variables that we keep during the analysis.

As in a standard abstract interpretation framework, our analysis carries out the fixed-point iteration until a fixed-point $(Pre_i, Post_i)$ (for some i) is reached. To infer the pre-conditions, our abstract semantics is equipped with bi-abduction over the combined domain. To allow the discovery of more precise preconditions, our abstraction procedure is also equipped with abduction, yielding the novel *abductive abstraction* (abs_a) for precondition discovery. The postcondition inference still employs the normal abstraction mechanism (abs).¹

We first set the precondition as `emp` and postcondition as `false` which signifies that we know nothing about the method (line 1). Then for each iteration, a forward bi-abductive analysis is employed to compute a new pre-/post-condition (line 4) based on the current specification. The analysis performs abstraction on both pre-/post-conditions obtained to maintain the finiteness of the shape domain. The obtained results are joined with the results from the previous iteration (line 6), and a widening is conducted over both to ensure termination of the analysis (line 7). If the analysis cannot continue due to a program bug, or cannot keep the number of shared logical variables/cutpoints (counted by `cp_no`) within a specified bound (n), then a failure is reported (line 8). At the end of each iteration, the inferred summary is used to update the specification of mn (line 9), which will be used for recursive calls (if any) of mn in next iteration. Finally we judge whether a fixed-point is already reached (line 10). The last few lines (from line 11) ensure that inferred specifications are indeed sound using a standard abstract semantics (without abduction). Any unsound specifications will be ruled out.

Intuitively, the $join^\dagger$ operator is applied over two abstract states, and tries to find a common shape as an abstraction for the separation part of both states. If such common shape is found, it performs convex hull and bag join for the pure parts. Otherwise it keeps a disjunction of the two states. The $widen^\dagger$ is analogous to $join^\dagger$. The difference is that we expect the heap portion of the first state is subsumed by the second one, and

¹ The analysis uses lifted versions of these operations (indicated by \dagger), which will be explained in more detail later.

```

Fixpoint Computation in the Combined Domain
Input:  $\mathcal{T}$ ,  $t\ mn\ ((t\ x)^*; (t\ y)^*)\ \{e\}$ ,  $n$ 
Local:  $i := 0$ ;  $Pre_i := \mathbf{emp}$ ,  $Post_i := \mathbf{false}$ ;
1   $\mathcal{T}' := \mathcal{T} \cup \{t\ mn\ ((t\ x)^*; (t\ y)^*)\ \text{requires } Pre_0\ \text{ensures } Post_0\ \{e\}\}$ ;
2  repeat
3     $i := i + 1$ ;
4     $(Pre_i, Post_i) := \llbracket e \rrbracket_{\mathcal{T}'}^A(Pre_{i-1}, Post_{i-1})$ ;
5     $(Pre_i, Post_i) := (\mathbf{abs}_a^\dagger(Pre_i), \mathbf{abs}^\dagger(Post_i))$ ;
6     $(Pre_i, Post_i) := (\mathbf{join}^\dagger(Pre_{i-1}, Pre_i), \mathbf{join}^\dagger(Post_{i-1}, Post_i))$ ;
7     $(Pre_i, Post_i) := (\mathbf{widen}^\dagger(Pre_{i-1}, Pre_i), \mathbf{widen}^\dagger(Post_{i-1}, Post_i))$ ;
8    if  $Pre_i = \mathbf{false}$  or  $Post_i = \mathbf{false}$  or  $\mathbf{cp\_no}(Pre_i) > n$  or  $\mathbf{cp\_no}(Post_i) > n$ 
then return fail end if
9     $\mathcal{T}' := \mathcal{T} \cup \{t\ mn\ ((t\ x)^*; (t\ y)^*)\ \text{requires } Pre_i\ \text{ensures } Post_i\ \{e\}\}$ ;
10 until  $\mathcal{T}'$  does not change
11  $Post = \llbracket e \rrbracket_{\mathcal{T}'} Pre_i$ ;
12 if  $Post = \mathbf{false}$  or  $Post \not\vdash Post_i * \mathbf{true}$  then return fail
13 else return  $\mathcal{T}'$ 
14 end if

```

Fig. 5. Main analysis algorithm.

then it applies the pure widening for the pure part. The formal definitions of \mathbf{join}^\dagger and \mathbf{widen}^\dagger and other details are left in our report [13] due to page limit.

Bi-Abductive Abstract Semantics. As shown in Fig. 5, our analysis employs two abstract semantics: a bi-abductive abstract semantics (i.e. the one equipped with abduction) (line 4), and an underlying abstract semantics (i.e. the one without abduction) (line 11). We first give the definition of the underlying semantics. Its type is defined as

$$\llbracket e \rrbracket : \text{AllSpec} \rightarrow \mathcal{P}_{\text{SH}} \rightarrow \mathcal{P}_{\text{SH}}$$

where AllSpec contains procedure specifications. For some program e and its given precondition Δ , the semantics calculates the postcondition $\llbracket e \rrbracket_{\mathcal{T}} \Delta$, for a given set of method specifications \mathcal{T} .

The essential constituents of the underlying semantics are the basic transition functions from a conjunctive abstract state (σ) to a conjunctive or disjunctive abstract state (σ or Δ) below:

$$\begin{array}{ll}
\mathbf{unfold}(x) : \text{SH} \rightarrow \mathcal{P}_{\text{SH}[x]} & \text{Unfolding} \\
\mathbf{exec}(d[x]) : \text{AllSpec} \rightarrow \text{SH}[x] \rightarrow \mathcal{P}_{\text{SH}} & \text{Heap-sensitive execution} \\
\mathbf{exec}(d) : \text{AllSpec} \rightarrow \text{SH} \rightarrow \mathcal{P}_{\text{SH}} & \text{Heap-insensitive execution}
\end{array}$$

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 ($x \mapsto c(v^*)$), and $\mathcal{P}_{\text{SH}[x]}$ contains all the (disjunctive) abstract states, each of which is composed by such conjunctive states. Here $\mathbf{unfold}(x)$ rearranges 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 $\mathbf{exec}(d[x])$. The third function defined for other (heap insensitive) commands d does not require such exposure of x .

The unfolding function is defined by the following two rules:

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

The symbolic execution of heap-sensitive commands $d[x]$ (i.e. $x.f_i$, $x.f_i := w$, or $\text{free}(x)$) assumes that the rearrangement $\text{unfold}(x)$ has been done prior to execution:

$$\frac{\sigma \vdash x \mapsto c(v_1, \dots, v_n) * \sigma'}{\text{exec}(x.f_i)(\mathcal{T})\sigma \rightsquigarrow \sigma \wedge \text{res}=v_i} \quad \frac{\sigma \vdash x \mapsto c(u^*) * \sigma'}{\text{exec}(\text{free}(x))(\mathcal{T})\sigma \rightsquigarrow \sigma'}$$

$$\frac{\sigma \vdash x \mapsto c(v_1, \dots, v_n) * \sigma'}{\text{exec}(x.f_i := w)(\mathcal{T})\sigma \rightsquigarrow \sigma' * x \mapsto c(v_1, \dots, v_{i-1}, w, v_{i+1}, \dots, v_n)}$$

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

$$\text{exec}(k)(\mathcal{T})\sigma =_{df} \sigma \wedge \text{res}=k \quad \frac{\text{isdatat}(c)}{\text{exec}(\text{new } c(v^*))(\mathcal{T})\sigma =_{df} \sigma * c(\text{res}, v^*)}$$

$$\frac{t \text{ mn } ((t_i u_i)_{i=1}^m; (t'_i v_i)_{i=1}^n) \text{ requires } \Phi_{pr} \text{ ensures } \Phi_{po} \in \mathcal{T} \quad \rho = [x'_i/u_i]_{i=1}^m \circ [y'_i/v_i]_{i=1}^n \quad \sigma \vdash \rho \Phi_{pr} * \sigma' \quad \rho_o = [r_i/v_i]_{i=1}^n \circ [x'_i/u'_i]_{i=1}^m \circ [y'_i/v'_i]_{i=1}^n \quad \rho_l = [r_i/y'_i]_{i=1}^n \text{ fresh logical } r_i}{\text{exec}(\text{mn}(x_1, \dots, x_m; y_1, \dots, y_n))(\mathcal{T})\sigma \rightsquigarrow (\rho_l \sigma') * (\rho_o \Phi_{po})}$$

The first three rules deal with constant (k), variable (x) and data node creation ($\text{new } c(v^*)$), respectively, while the last rule handles method invocation. The test $\text{isdatat}(c)$ returns true iff c is a data node. In the last rule, the call site is ensured to meet the precondition of mn , as signified by $\sigma \vdash \rho \Phi_{pr} * \sigma'$. In this case, the execution succeeds and the post-state of the method call involves mn 's postcondition as signified by $\rho_o \Phi_{po}$.

A lifting function \dagger is defined to lift unfold 's domain to \mathcal{P}_{SH} :

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

and this function is overloaded for exec to lift both its domain and range to \mathcal{P}_{SH} :

$$\text{exec}^\dagger(d)(\mathcal{T}) \bigvee \sigma_i =_{df} \bigvee (\text{exec}(d)(\mathcal{T})\sigma_i)$$

Based on the transition functions above, we can define the abstract semantics for a program e as follows (where loops are already translated into tail-recursions):

$$\begin{aligned} \llbracket d[x] \rrbracket_{\mathcal{T}} \Delta &=_{df} \text{exec}^\dagger(d[x])(\mathcal{T}) \circ \text{unfold}^\dagger(x) \Delta \\ \llbracket d \rrbracket_{\mathcal{T}} \Delta &=_{df} \text{exec}^\dagger(d)(\mathcal{T}) \Delta \\ \llbracket e_1; e_2 \rrbracket_{\mathcal{T}} \Delta &=_{df} \llbracket e_2 \rrbracket_{\mathcal{T}} \circ \llbracket e_1 \rrbracket_{\mathcal{T}} \Delta \\ \llbracket x := e \rrbracket_{\mathcal{T}} \Delta &=_{df} [x'/x, r'/\text{res}](\llbracket e \rrbracket_{\mathcal{T}} \Delta) \wedge x=r' \text{ fresh logical } x', r' \\ \llbracket \text{if } (v) e_1 \text{ else } e_2 \rrbracket_{\mathcal{T}} \Delta &=_{df} (\llbracket e_1 \rrbracket_{\mathcal{T}}(v \wedge \Delta)) \vee (\llbracket e_2 \rrbracket_{\mathcal{T}}(\neg v \wedge \Delta)) \end{aligned}$$

We shall now present the definitions of our bi-abductive abstract semantics. Its type is

$$\llbracket e \rrbracket^A : \text{AllSpec} \rightarrow (\mathcal{P}_{\text{SH}} \times \mathcal{P}_{\text{SH}}) \rightarrow (\mathcal{P}_{\text{SH}} \times \mathcal{P}_{\text{SH}})$$

It takes a piece of program code and a specification table, and maps a pair of (disjunctive) set of symbolic heaps to another such pair (where the first in the pair is the accumulated precondition and the second is the current state). It relies on the following two basic functions:

$$\begin{aligned} \text{Unfold}(x) &: (\text{SH} \times \text{SH}) \rightarrow (\mathcal{P}_{\text{SH}} \times \mathcal{P}_{\text{SH}}) \\ \text{Exec}(ds) &: \text{AllSpec} \rightarrow (\text{SH} \times \text{SH}) \rightarrow (\mathcal{P}_{\text{SH}} \times \mathcal{P}_{\text{SH}}) \end{aligned}$$

The definitions of both functions are given below:

$$\begin{aligned}
& \text{Unfold}(x)(\sigma', \sigma) =_{df} \\
& \quad \text{if } (\sigma * [\sigma_m] \triangleright x \mapsto c(y^*) * \text{true for fresh logical vars } y^*) \wedge (\sigma' * \sigma_m \not\vdash \text{false}) \\
& \quad \text{then let } \Delta = \text{unfold}(x)(\sigma * \sigma_m) \text{ in } (\sigma' * \sigma_m, \Delta) \\
& \quad \text{else } (\sigma', \text{false}) \\
& \text{Exec}(ds)(\mathcal{T})(\sigma', \sigma) =_{df} \text{let } \Delta = \text{exec}(ds)(\mathcal{T})\sigma \text{ in } (\sigma', \Delta) \\
& \quad \text{where } ds \text{ is either } d[x] \text{ or } d, \text{ except for procedure call} \\
& \quad \frac{t \text{ mn } ((t_i u_i)_{i=1}^m; (t'_i v_i)_{i=1}^n) \text{ requires } \Phi_{pr} \text{ ensures } \Phi_{po} \in \mathcal{T} \\
& \quad \quad \rho = [x'_i/u_i]_{i=1}^m \circ [y'_i/v_i]_{i=1}^n \circ \sigma * [\sigma'_1] \triangleright \rho \Phi_{pr} * \sigma_1 \quad \sigma' * \sigma'_1 \not\vdash \text{false} \\
& \quad \quad \rho_o = [r_i/v_i]_{i=1}^n \circ [x'_i/u_i]_{i=1}^m \circ [y'_i/v'_i]_{i=1}^n \quad \rho_l = [r_i/y'_i]_{i=1}^n \quad \text{fresh logical vars } r_i}{\text{Exec}(\text{mn}(x_{1..m}; y_{1..n}))(\mathcal{T})(\sigma', \sigma) =_{df} (\sigma' * \sigma'_1, (\rho_o \Phi_{po}) * (\rho_l \sigma_1))}
\end{aligned}$$

The Unfold function firstly tests (using bi-abduction) whether the node $x \mapsto c(y^*)$ is in σ , if not, abduction is applied to find the missing σ_m . If σ' and σ_m do not contradict, it unfolds $\sigma * \sigma_m$ to expose x (via the unfold function defined earlier in this section), and adds σ_m to precondition. Otherwise, it returns false for the current state.

The Exec function symbolically executes the command ds (via the exec function defined earlier in this section) and translates the current state σ to a disjunction of new states Δ . The special case is the method invocation, which may require bi-abduction to be applied for the current state. When the method mn is invoked, we take its current specification (Φ_{pr}, Φ_{po}) from \mathcal{T} , and substitute the formal parameters u_i and v_i by the current arguments x'_i and y'_i respectively. Note that prime notations x'_i and y'_i denote the current values of x_i and y_i in the current state σ . Then we apply bi-abduction from the current state σ to the precondition $\rho \Phi_{pr}$. If it succeeds, the discovered missing state σ'_1 will be propagated back to the precondition σ' to help make the symbolic execution to succeed. The postcondition of mn , Φ_{po} is substituted by ρ_o in order to be added to the current state. Since the variables y_i are call-by-reference, we let r_i to be the intermediate variables, while the variables y'_i denote the latest values.

A lifting function \dagger is defined to lift Unfold's and Exec's domains:

$$\begin{aligned}
& \text{Unfold}^\dagger(x) \vee (\sigma'_i, \sigma_i) =_{df} \vee (\text{Unfold}(x)(\sigma'_i, \sigma_i)) \\
& \text{Exec}^\dagger(ds)(\mathcal{T}) \vee (\sigma'_i, \sigma_i) =_{df} \vee (\text{Exec}(ds)(\mathcal{T})(\sigma'_i, \sigma_i))
\end{aligned}$$

Based on the above functions, the bi-abductive abstract semantics is defined as follows:

$$\begin{aligned}
& \llbracket d[x] \rrbracket_{\mathcal{T}}^\Delta(\Delta', \Delta) =_{df} \text{Exec}^\dagger(d[x])(\mathcal{T}) \circ \text{Unfold}^\dagger(x)(\Delta', \Delta) \\
& \llbracket d \rrbracket_{\mathcal{T}}^\Delta(\Delta', \Delta) =_{df} \text{Exec}^\dagger(d)(\mathcal{T})(\Delta', \Delta) \\
& \llbracket e_1; e_2 \rrbracket_{\mathcal{T}}^\Delta(\Delta', \Delta) =_{df} \llbracket e_2 \rrbracket_{\mathcal{T}}^\Delta \circ \llbracket e_1 \rrbracket_{\mathcal{T}}^\Delta(\Delta', \Delta) \\
& \llbracket x := e \rrbracket_{\mathcal{T}}^\Delta(\Delta', \Delta) =_{df} [x'/x, r'/\text{res}](\llbracket e \rrbracket_{\mathcal{T}}^\Delta(\Delta', \Delta \wedge x=r')) \quad \text{fresh logical } x', r' \\
& \llbracket \text{if } (v) e_1 \text{ else } e_2 \rrbracket_{\mathcal{T}}^\Delta(\Delta', \Delta) =_{df} (\llbracket e_1 \rrbracket_{\mathcal{T}}^\Delta(\Delta', v \wedge \Delta)) \vee (\llbracket e_2 \rrbracket_{\mathcal{T}}^\Delta(\Delta', \neg v \wedge \Delta))
\end{aligned}$$

Abductive Abstraction. As we mentioned earlier in the `merge` example, to verify such programs may require very precise preconditions that a standard abstraction mechanism may fail to achieve. To cater for such a need, we design a novel *abductive abstraction* function abs_a , which equips abstraction with an abductive reasoning capacity where necessary. In such scenarios, user-specified predicates can offer some guidance in the

abstraction in order to discover extra data structure properties for precondition. The new abductive abstraction function is given as follows:

$$\begin{array}{l}
\text{abs}_a(\sigma \wedge x_0=e) =_{df} \sigma[e/x_0] \\
\text{abs}_a(\sigma \wedge e=x_0) =_{df} \sigma[e/x_0] \\
\frac{x_0 \notin \text{Reach}(\sigma)}{\text{abs}_a(\mathbb{H}(c)(x_0, v^*) * \sigma) =_{df} \sigma * \text{true}} \\
\frac{p_2(u_2^*) \equiv \Phi \quad \mathbb{H}(c_1)(x, v_1^*) * \sigma_1 \vdash p_2(x, v_2^*) \wedge \pi_2 \quad \text{Reach}(p_2(x, v_2^*) \wedge \pi_2 * \sigma_3) \cap \{v_1^*\} = \emptyset}{\text{abs}_a(\mathbb{H}(c_1)(x, v_1^*) * \sigma_1 * \sigma_3) =_{df} p_2(x, v_2^*) \wedge \pi_2 * \sigma_3} \\
\frac{\mathbb{H}(c_1)(x, v_1^*) * \sigma_1 * [\sigma'] \triangleright p_2(x, v_2^*) \wedge \pi_2 \quad \text{Reach}(p_2(x, v_2^*) \wedge \pi_2 * \sigma_3) \cap \{v_1^*\} = \emptyset}{\text{abs}_a(\mathbb{H}(c_1)(x, v_1^*) * \sigma_1 * \sigma_3) =_{df} p_2(x, v_2^*) \wedge \pi_2 * \sigma_3}
\end{array}$$

where $\mathbb{H}(c)(x, v^*)$ denotes $x \mapsto c(v^*)$ if c is a data node or $c(x, v^*)$ if c is a predicate. The function $\text{Reach}(\sigma)$ returns all pointer variables which are reachable from free variables in the abstract state σ . The first two rules eliminate logical variables, and the third rule drops heap garbage that is unreachable from program variables. The fourth rule combines shape formulae and eliminate logical pointer variables which are not reachable from other program variables. The predicate p_2 is selected from the user-defined predicates environments and it is the target shape to be abstracted to.

The last rule applies when the state $\mathbb{H}(c_1)(x, v_1^*) * \sigma_1$ cannot be abstracted to the predicate p_2 using standard abstraction but can be abstracted to predicate p_2 with the help of abductive reasoning. When applying such an abstraction function during the precondition discovery, the extra information σ' discovered by abduction will be propagated back to the precondition to improve the precision.

The lifting function is applied for abs_a to lift both its domain and range to disjunctive abstract states $\mathcal{P}_{\text{SH}}: \text{abs}_a^\dagger \vee \sigma_i =_{df} \vee \text{abs}_a(\sigma_i)$, allowing it to be used in the analysis.

The soundness and termination of our analysis are given in the technical report [13].

6 Experiments and Evaluation

We have implemented a prototype system and evaluated it over a number of heap-manipulating programs to test the viability and precision of our approach. Our experimental results were achieved with an Intel Core 2 Quad CPU 2.66GHz with 8GB RAM. We have also defined a library of predicates covering popular data structures and variety of properties. These properties can be grouped in the following categories: MS (*memory safety*): all memory accesses are safe, no dangling/null pointers dereferences; SC (*same content*): the content of the final data structure remains the same as that of the input data structure; IN (*insertion*): the input data is inserted into the final data structure; SO (*sorted*): data structures are sorted according to a criterion, eg. in case of a list each node's content is less than or equal to its successor's; BS (*binary search*): data structures are binary search trees; DL (*double-linked list*): data structures are double-linked lists; and AL (*AVL tree*): data structures are AVL trees. The predicates required as input by our tool can be selected from the library or can be supplied by users, according to the input program data structures and the properties of interest. Usually, the upper bound of cutpoints is set to be twice the number of input program variables to improve the precision. Some of our results are presented in Table 1.

Prog.	LOC	Time	Prop	Prog.	LOC	Time	Prop
Singly Linked List				Doubly Linked List			
create	10	1.12	MS	create	15	1.47	MS/DL
delete	9	1.20	MS/SO	append	24	2.53	MS/DL/SC/SO
insert	9	1.16	MS/SO/IN	insert	22	2.32	MS/DL/IN/SO
traverse	9	1.35	MS/SO/SC	Binary Search Tree			
length	11	1.28	MS/SO/SC	create	18	2.58	MS/BS
append	11	1.47	MS/SO/SC	delete	48	4.76	MS/BS
take	12	1.28	MS/SO/SC	insert	22	3.57	MS/BS/IN
reverse	13	1.72	MS/SC	search	22	2.78	MS/BS/SC
filter	15	2.37	MS/SO	height	15	1.56	MS/BS/SC
Sorting algorithm				count	17	1.63	MS/BS/SC
insert_sort	32	2.72	MS/SC/SO	flatten	32	2.74	MS/BS/DL/SC/SO
merge_sort	78	4.18	MS/SC/SO	AVL Tree			
quick_sort	70	5.72	MS/SC/SO	insert	114	27.57	MS/BS/AL/IN
select_sort	45	3.16	MS/SC/SO	delete	239	34.42	MS/BS/AL

Table 1. Experimental Results. The column **LOC** is for the number of program lines; **Time** expresses our tool running time (in seconds); **Prop** denotes the inferred specification properties.

In comparison to previous approaches, the first observation concerns the precision of our analysis. Since our tool uses a combined domain it can discover more expressive specifications to guarantee memory safety and functional correctness. For example in case of the `take` program which traverses the list down for a user-specified number `n` of nodes, we can find that the input list length must be no less than `n`. However the previous tools based on shape domains (like Abductor [5]) can only discover a precondition that requires the input list to be non-empty which would not be sufficient to guarantee memory safety. Moreover more complex functional properties regarding the data structures content (like `SO` for merge program but in general for all sorting programs) can also not be discovered by the previous tools (like Abductor) based on a simple shape domain. There are other tools (like Xisa [6] or Thor [18]) that can work on a combined domain but require certain annotations to guide their analysis. Thor [18] requires shape information for each input parameter and Xisa [6] requires shape information for program variables used in loops. Since our shape domain includes tree data structures, our tool is able to discover complex functional specifications for binary search trees and AVL trees in contrast to the previous approaches. For example in case of the `flatten` program our tool is able to discover that the input data structure is a binary search tree while the output data structure is a sorted doubly linked list having the same data content (values stored inside the nodes) as that of the input.

The second observation regarding our experimental results is that the analysis may discover more than one correct specification for some programs. For example, given two predicates, ordinary linked list and sorted list, we can obtain two specifications for most of the sorting algorithms. When there are more than one user-supplied predicate definitions, the analysis can have multiple choices during the abstraction. Multiple specifications can be useful in program verification, e.g. the sorted version for the `append` method, where the two input lists and the output list are all sorted, is useful in the verification of `quick_sort`, while the sorted list version for the `insert` method is also useful to help verify the functional correctness of `insert_sort`.

7 Related Work and Conclusion

Dramatic advances have been made in synthesising specifications for heap-manipulating programs. The local shape analysis [8] infers loop invariants for list-processing programs, followed by the SpaceInvader/Abductor tool to infer full method specifications over the separation domain, so as to verify pointer safety for larger industrial codes [5, 26]. The SLayer tool [9] implements an inter-procedural analysis for programs with shape information. A combination of shape and bag abstraction is used in [25] to verify linearizability. Compared with them, our abstraction is more general since it is driven by predicates and is not restricted to linked lists. To deal with size information (such as number of nodes in lists/trees), Thor [18] transfers a heap-processing program to a numerical one, so that size properties can be obtained by further analysis. A similar approach [10] combines a set domain (for shape) with its cardinality domain (for corresponding numerical information) in a more general framework. Compared with these works, our approach can discover specifications with stronger invariants such as sortedness and bag-related properties, which have not been addressed in the previous works. The analyses [6, 19, 20] can all handle shape and numerical information over a combined domain, but require user given preconditions for the program whereas here we compute the whole specification at once. Recently, Rival and Chang [23] propose an inductive predicate to summarise call stacks along with heap structures in a context of a whole-program analysis. In contrast our analysis is modular.

There are also other approaches that can synthesise shape-related program invariants. The shape analysis framework TVLA [24] is based on three-valued logic. It is capable of handling complicated data structures and properties, such as sortedness. Guo et al. [11] report a global shape analysis that discovers inductive structural shape invariants from the code. Kuncak et al. [15] develop a role system to express and track referencing relationships among objects. Hackett and Rugina [12] can deal with AVL-trees but is customised to handle only tree-like structures with height property. Bouajjani et al. [2, 3] propose a program analysis in an abstract domain with SL3 (Singly-Linked List Logic) and size, sortedness and multi-set properties. However, their heap domain is restricted to singly-linked list only, and their shape analysis is separated from numerical and multi-set analyses. Compared with these works, separation logic based approaches benefit from the frame rule with support for local reasoning.

There are also approaches which unify reasoning over shape and data using either a combination of appropriate decision procedures inside Satisfiability-Modulo-Theories (SMT) solvers (e.g. [21, 16]) or a combination of appropriate abstract interpreters inside a software model checker (e.g. [1]). Compared with our work, their heap domains are mainly restricted to linked lists.

Conclusion. We have reported a program analysis which automatically discovers program specifications over a combined separation and pure(numerical and bag) domain. The novel components of our analysis include an abductive abstract semantics and an abductive abstraction mechanism (for precondition discovery) in the combined domain. We have built a prototype system and the initial experimental results are encouraging.

Acknowledgement. This work was supported in part by EPSRC project EP/G042322.

References

1. Beyer, D., Henzinger, T.A., Théoduloz, G.: Configurable software verification: Concretizing the convergence of model checking and program analysis. In: CAV (2007)
2. Bouajjani, A., Dragoi, C., Enea, C., Sighireanu, M.: On inter-procedural analysis of programs with lists and data. In: PLDI (2011)
3. Bouajjani, A., Dragoi, C., Enea, C., Sighireanu, M.: Abstract domains for automated reasoning about list-manipulating programs with infinite data. In: VMCAI (2012)
4. Bozga, M., Iosif, R., Lakhnech, Y.: Storeless semantics and alias logic. In: PEPM (2003)
5. Calcagno, C., Distefano, D., O'Hearn, P., Yang, H.: Compositional shape analysis by means of bi-abduction. *J. ACM* 58(6) (2011)
6. Chang, B.Y.E., Rival, X.: Relational inductive shape analysis. In: POPL (2008)
7. Chin, W.N., David, C., Nguyen, H.H., Qin, S.: Automated verification of shape, size and bag properties via user-defined predicates in separation logic. *Sci. of Comp. Prog.* 77 (2012)
8. Distefano, D., O'Hearn, P.W., Yang, H.: A local shape analysis based on separation logic. In: TACAS (2006)
9. Gotsman, A., Berdine, J., Cook, B.: Interprocedural shape analysis with separated heap abstractions. In: SAS (2006)
10. Gulwani, S., Lev-Ami, T., Sagiv, M.: A combination framework for tracking partition sizes. In: Shao, Z., Pierce, B.C. (eds.) POPL (2009)
11. Guo, B., Vachharajani, N., August, D.I.: Shape analysis with inductive recursion synthesis. In: PLDI (2007)
12. Hackett, B., Rugina, R.: Region-based shape analysis with tracked locations. In: POPL (2005)
13. He, G., Qin, S., Chin, W.N., Craciun, F.: Automated specification discovery in a combined abstract domain - reseach report. (2012), <http://pls.tees.ac.uk/~guan/fullspec/techreport.pdf>
14. Jonkers, H.: Abstract storage structures. In: Algorithmic Languages (1981)
15. Kuncak, V., Lam, P., Rinard, M.C.: Role analysis. In: POPL (2002)
16. Lahiri, S.K., Qadeer, S.: Back to the future: revisiting precise program verification using smt solvers. In: POPL (2008)
17. Magill, S., Tsai, M.H., Lee, P., Tsay, Y.K.: Thor: A tool for reasoning about shape and arithmetic. In: CAV (2008)
18. Magill, S., Tsai, M.H., Lee, P., Tsay, Y.K.: Automatic numeric abstractions for heap-manipulating programs. In: POPL (2010)
19. Qin, S., He, G., Luo, C., Chin, W.N., Chen, X.: Loop invariant synthesis in a combined abstract domain. *Journal of Symbolic Computation* 50 (2013)
20. Qin, S., Luo, C., Chin, W.N., He, G.: Automatically refining partial specifications for program verification. In: FM (2011)
21. Rakamaric, Z., Bruttomesso, R., Hu, A.J., Cimatti, A.: Verifying heap-manipulating programs in an smt framework. In: ATVA (2007)
22. Reynolds, J.C.: Separation logic: A logic for shared mutable data structures. In: LICS (2002)
23. Rival, X., Chang, B.Y.E.: Calling context abstraction with shapes. In: POPL (2011)
24. Sagiv, M., Reps, T.W., Wilhelm, R.: Parametric shape analysis via 3-valued logic. *ACM Trans. Program. Lang. Syst.* 24(3) (2002)
25. Vafeiadis, V.: Shape-value abstraction for verifying linearizability. In: VMCAI (2009)
26. Yang, H., Lee, O., Berdine, J., Calcagno, C., Cook, B., Distefano, D., O'Hearn, P.W.: Scalable shape analysis for systems code. In: CAV (2008)