

# Loop Invariant Synthesis in a Combined Domain

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

<sup>1</sup> Teesside University, Middlesbrough, TS1 3BA, UK

<sup>2</sup> Durham University, Durham, DH1 3LE, UK

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

**Abstract.** Automated verification of memory safety and functional correctness for heap-manipulating programs has been a challenging task, especially when dealing with complex data structures with strong invariants involving both shape and numerical properties. Existing verification systems usually rely on users to supply annotations, which can be tedious and error-prone and can significantly restrict the scalability of the verification system. In this paper, we reduce the need of user annotations by automatically inferring loop invariants over an abstract domain with both separation and numerical information. Our loop invariant synthesis is conducted automatically by a fixpoint iteration process, equipped with newly designed abstraction mechanism, and join and widening operators. Initial experiments have confirmed that we can synthesise loop invariants with non-trivial constraints.

## 1 Introduction

Although research on software verification has a long and distinguished history dating back to the 1960's, it remains a challenging problem to automatically verify heap manipulating programs written in mainstream imperative languages. This is in part due to the shared mutable data structures lying in programs, and the need to track various program properties, such as structural numerical information (length and height) and relational numerical information (sortedness and binary search tree properties).

Since the emergence of separation logic [14, 25], dramatic advances have been made in automated software verification, e.g. the Smallfoot tool [1] for the verification on pointer safety (i.e. shape properties asserting that pointers cannot go wrong), the verification on termination [2], the verification for object-oriented programs [6, 22], and Dafny [18] and HIP/SLEEK [5, 20, 21] for more general properties (both structural and numerical ones) for heap-manipulating programs.

These verification systems generally require users to provide specifications for each method as well as invariants for each loop, which is both tedious and error-prone. This also affects their scalability, as there can be many methods in a program and each method may still contain several while loops.

To conquer this problem, separation logic based shape analysis techniques are brought in, e.g., the SpaceInvader tool [3, 9, 27]. As a further step of Smallfoot, it automatically infers method specifications and loop invariants for pointer

---

\* Now with Citigroup Inc.

safety in the shape domain. Other works such as THOR [19] incorporate simple numerical information into the shape domain to allow automated synthesis of properties like list length. Their success proves the necessity and feasibility for shape analysis to help automate the verification process.

However, the prior analyses focus mainly on relatively simple properties, such as pointer safety for lists and list length information. It is difficult to apply them in the presence of more sophisticated program properties, such as:

- More flexible user-defined data structures, such as trees;
- Relational numerical properties, like sortedness and binary search property.

These properties can be part of the full functional correctness of heap-manipulating programs. The (aforementioned) HIP/SLEEK tool aims to verify such properties and it allows users to define their own shape predicates to express their desired level of correctness.

In this paper, we make the first stride to improve the level of automation for HIP/SLEEK-like verification systems by discovering loop invariants automatically over the combined shape and numerical domain. This proves to be a challenging problem especially since we aim towards full functional correctness that HIP/SLEEK targets at. Our approach is based on the framework of abstract interpretation [7] with fixpoint computation. It makes the following contributions in summary:

- We propose a loop invariant synthesis with novel operations for abstraction, join and widening over a combined shape and numerical domain.
- We demonstrate that our analysis is sound w.r.t. concrete program semantics and always terminates.
- We have integrated our solution with HIP/SLEEK and conducted some initial experiments. The experimental results confirm the viability of our solution and show that we can effectively eliminate the need for user-provision of loop invariants which were previously necessary in verification.

We shall next illustrate our approach informally via an example before presenting the formal details.

## 2 The Approach

Before giving an illustrative example for the analysis, we will first introduce our specification mechanism which follows the HIP/SLEEK system.

### 2.1 Separation Logic and User-defined Predicates

Separation logic [14, 25] extends Hoare logic to support reasoning about shared mutable data structures. It adds two more connectives to classical logic: separation conjunction  $*$  and spatial implication  $-*$ . The formula  $\mathbf{p}_1 * \mathbf{p}_2$  asserts that two heaps described by the formulae  $\mathbf{p}_1$  and  $\mathbf{p}_2$  are domain-disjoint, while  $\mathbf{p}_1 -* \mathbf{p}_2$  asserts that if the current heap is extended with a disjoint heap described by

the formula  $p_1$ , then the formula  $p_2$  holds in the extended heap. In this paper we only use separation conjunction.

Similar to the HIP/SLEEK system, we allow user-defined inductive predicates to specify both separation and numerical properties. For example, with a data structure definition for a node in a list `data node { int val; node next; }`, we can define a predicate for a list as

$$\text{root}::\text{ll}\langle n \rangle \equiv (\text{root}=\text{null} \wedge n=0) \vee (\exists v, q, m \cdot \text{root}::\text{node}\langle v, q \rangle * \text{q}::\text{ll}\langle m \rangle \wedge n=m+1)$$

The parameter `root` for the predicate `ll` is the root pointer referring to the list. Its length is denoted by `n`. A uniform notation  $p::c\langle v_1, \dots, v_k \rangle$  is used for either a singleton heap or a predicate. If `c` is a data node with fields  $f_1, \dots, f_k$ , the notation represents a singleton heap,  $p \mapsto c[f_1:v_1, \dots, f_k:v_k]$ , e.g. the `root::node` above. If `c` is a predicate name, then the data structure pointed to by `p` has the shape `c` with parameters  $v_1, \dots, v_k$ , e.g., the `q::ll` above.

We can also define a list segment as follows:

$$\text{ls}\langle p, n \rangle \equiv (\text{root}=p \wedge n=0) \vee (\text{root}::\text{node}\langle \_, q \rangle * \text{q}::\text{ls}\langle p, m \rangle \wedge n=m+1)$$

where we use the following shortened notation: (i) default `root` parameter in LHS may be omitted, (ii) unbound variables, such as `q` and `m`, are implicitly existentially quantified, and (iii) the underscore `_` denotes an existentially quantified anonymous variable.

If the user wants to verify a sorting algorithm, they can incorporate sortedness property into the above predicates as follows:

$$\begin{aligned} \text{sll}\langle n, mn, mx \rangle &\equiv (\text{root}::\text{node}\langle mn, \text{null} \rangle \wedge n=1 \wedge mn=mx) \vee \\ &\quad (\text{root}::\text{node}\langle mn, q \rangle * \text{q}::\text{sll}\langle n_1, k, mx \rangle \wedge mn \leq k \wedge n=n_1+1) \\ \text{sls}\langle p, n, mn, mx \rangle &\equiv (\text{root}::\text{node}\langle mn, p \rangle \wedge n=1 \wedge mn=mx) \vee \\ &\quad (\text{root}::\text{node}\langle mn, q \rangle * \text{q}::\text{sls}\langle p, n_1, k, mx \rangle \wedge mn \leq k \wedge n=n_1+1) \end{aligned}$$

where `mn` and `mx` denote resp. the min and max values stored in the sorted list. Such user-supplied predicates can be used to specify loop invariants and method pre/post-specifications.

## 2.2 Illustrative Example

We now illustrate via an example our loop invariant synthesis process. The method `ins_sort` (Figure 1) sorts a linked list with the insertion sort algorithm. It is implemented with two nested while loops. The outer loop traverses the whole list `x`, takes out each node from it (line 7), and inserts that node into another already sorted list `r` (which is empty initially before the sorting). This insertion process makes use of the inner while loop in lines 9-11 to look for a proper position in the already sorted list for the new node to be inserted. The actual insertion takes place at lines 12-14.

To verify this program, we need to synthesise appropriate loop invariants for both while loops. Our analysis follows a standard fixpoint iteration process. It

<pre> 0 data node { int val;                 node next; } 1 node ins_sort(node x) 2   requires x::ll⟨n⟩ 3   ensures res::sll⟨n, mn, mx⟩ 4 {int v; 5   node r, cur, srt, prv=null; 6   while (x != null) { 7     cur=x; x=x.next; v=cur.val; 8     srt=r; prv=null; </pre>	<pre> 9   while (srt != null &amp;&amp;           srt.val &lt;= v) { 10    prv=srt; srt=srt.next; 11  } 12  cur.next=srt; 13  if (prv != null) prv.next=cur; 14  else r=cur; 15  } 16  return r; 17 } </pre>
---	--

**Fig. 1.** Insertion sort for linked list.

starts with the (abstract) program state immediately before the while loop (i.e., the initial state) and symbolically executes the loop body for several iterations, until the obtained states converge to a fixpoint, which is the loop invariant.<sup>1</sup> At the start of each iteration, the obtained state from the previous iteration is joined with the initial state. In addition to this join operator, we have also defined an abstraction function and a widening operator both of which will help the fixpoint iteration to converge. The join and widening operators are specifically designed to handle both shape and numerical information.

As for our example, due to the presence of nested loops, each iteration of the analysis for the outer loop actually infers a loop invariant for the inner loop. We shall now illustrate how we synthesise a loop invariant for the inner loop.

Suppose that in one iteration for the outer loop, the state at line 9 becomes  $r::sll\langle n_r, a, b \rangle * cur::node\langle v, x \rangle * x::ll\langle n_x \rangle \wedge srt=r \wedge prv=null \wedge n_r+n_x+1=n$

Note that since the inner loop does not mutate the heap part referred to by  $cur$  and  $x$  (i.e.,  $cur::node\langle v, x \rangle * x::ll\langle n_x \rangle$ ), we can ignore it during the invariant synthesis and add it back to the program state using the frame rule of separation logic [25]. Therefore, the initial state for loop invariant synthesis becomes

$$r::sll\langle n_r, a, b \rangle \wedge srt=r \wedge prv=null \wedge n_r+n_x+1=n \quad (1)$$

From this state, symbolically executing the loop body once yields the state:

$$r::node\langle a, srt \rangle * srt::sll\langle n_s, c_1, b \rangle \wedge prv=r \wedge a \leq c_1 \wedge a \leq v \wedge n_r+1=n-n_x \wedge n_s+1=n_r \quad (2)$$

which says that pointer  $srt$  moves towards the tail of the list for one node. Then we join it with the initial state (1) to obtain

$$\begin{aligned} & (r::sll\langle n_r, a, b \rangle \wedge srt=r \wedge prv=null \wedge n_r+n_x+1=n) \vee \\ & (r::node\langle a, srt \rangle * srt::sll\langle n_s, c_1, b \rangle \wedge \\ & \quad prv=r \wedge a \leq c_1 \wedge a \leq v \wedge n_r+1=n-n_x \wedge n_s+1=n_r) \end{aligned} \quad (3)$$

The second iteration over the loop body starts with (3) and exhibits (also) the case that  $srt$  runs two nodes towards tail, while  $prv$  goes one node. Its result is then joined with pre-state (1) to become the current state:

<sup>1</sup> The fixpoint iteration converges if one more iteration still yields the same result.

$$(3) \vee \mathbf{r}::\mathbf{node}\langle \mathbf{a}, \mathbf{prv} \rangle * \mathbf{prv}::\mathbf{node}\langle \mathbf{c}_1, \mathbf{srt} \rangle * \mathbf{srt}::\mathbf{sll}\langle \mathbf{n}_s, \mathbf{c}_2, \mathbf{b} \rangle \wedge \mathbf{a} \leq \mathbf{c}_1 \leq \mathbf{c}_2 \wedge \mathbf{c}_1 \leq \mathbf{v} \wedge \mathbf{n}_r + 1 = \mathbf{n} - \mathbf{n}_x \wedge \mathbf{n}_s + 2 = \mathbf{n}_r \quad (4)$$

Executing the loop body a third time returns a post-state where three nodes are passed by  $\mathbf{srt}$ , and two by  $\mathbf{prv}$ , as below:

$$(4) \vee \mathbf{r}::\mathbf{node}\langle \mathbf{a}, \mathbf{r}_0 \rangle * \mathbf{r}_0::\mathbf{node}\langle \mathbf{c}_1, \mathbf{prv} \rangle * \mathbf{prv}::\mathbf{node}\langle \mathbf{c}_2, \mathbf{srt} \rangle * \mathbf{srt}::\mathbf{sll}\langle \mathbf{n}_s, \mathbf{c}_3, \mathbf{b} \rangle \wedge \mathbf{a} \leq \mathbf{c}_1 \leq \mathbf{c}_2 \leq \mathbf{c}_3 \wedge \mathbf{c}_2 \leq \mathbf{v} \wedge \mathbf{n}_r + 1 = \mathbf{n} - \mathbf{n}_x \wedge \mathbf{n}_s + 3 = \mathbf{n}_r$$

where we have an auxiliary logical variable  $\mathbf{r}_0$ . Following this trend, it is predictable that every iteration hereafter will introduce an additional logical variable (referring to a list node). If we indulge in such increase in the subsequent iterations, the analysis will never terminate. Our abstraction process prevents this from happening by eliminating such logical variables as follows:

$$(4) \vee \mathbf{r}::\mathbf{sls}\langle \mathbf{prv}, \mathbf{n}_1, \mathbf{a}, \mathbf{c}_1 \rangle * \mathbf{prv}::\mathbf{node}\langle \mathbf{c}_2, \mathbf{srt} \rangle * \mathbf{srt}::\mathbf{sll}\langle \mathbf{n}_s, \mathbf{c}_3, \mathbf{b} \rangle \wedge \mathbf{a} \leq \mathbf{c}_1 \leq \mathbf{c}_2 \leq \mathbf{c}_3 \wedge \mathbf{c}_2 \leq \mathbf{v} \wedge \mathbf{n}_r + 1 = \mathbf{n} - \mathbf{n}_x \wedge \mathbf{n}_s + 3 = \mathbf{n}_r \wedge \mathbf{n}_1 = 2$$

Note that the heap part  $\mathbf{r}::\mathbf{node}\langle \mathbf{a}, \mathbf{r}_0 \rangle * \mathbf{r}_0::\mathbf{node}\langle \mathbf{c}_1, \mathbf{prv} \rangle$  is abstracted as a sorted list segment  $\mathbf{r}::\mathbf{sls}\langle \mathbf{prv}, \mathbf{n}_1, \mathbf{a}, \mathbf{c}_1 \rangle$  with  $\mathbf{n}_1$  denoting the length of the segment and  $\mathbf{n}_1 = 2$  added into the state. This abstraction process ensures that our analysis does not allow the shape to increase infinitely.

The fourth iteration responds with a post-state where four nodes are passed by  $\mathbf{srt}$ , and three by  $\mathbf{prv}$ . Therefore an abstraction is performed to remove the logical pointer variables. As a simplification of the presentation, we denote  $\sigma$  as  $\mathbf{r}::\mathbf{sls}\langle \mathbf{prv}, \mathbf{n}_1, \mathbf{a}, \mathbf{c}_1 \rangle * \mathbf{prv}::\mathbf{node}\langle \mathbf{c}_2, \mathbf{srt} \rangle * \mathbf{srt}::\mathbf{sll}\langle \mathbf{n}_s, \mathbf{c}_3, \mathbf{b} \rangle \wedge \mathbf{a} \leq \mathbf{c}_1 \leq \mathbf{c}_2 \leq \mathbf{c}_3 \wedge \mathbf{c}_2 \leq \mathbf{v} \wedge \mathbf{n}_r + 1 = \mathbf{n} - \mathbf{n}_x$ , and the abstracted result (after the fourth iteration) is

$$(4) \vee (\sigma \wedge \mathbf{n}_s + 3 = \mathbf{n}_r \wedge \mathbf{n}_1 = 2) \vee (\sigma \wedge \mathbf{n}_s + 4 = \mathbf{n}_r \wedge \mathbf{n}_1 = 3)$$

for which we have an observation that the last two disjunctions share the same shape part (as in  $\sigma$ ). This disjunction will be transferred to the numerical domain, as follows:

$$(4) \vee (\sigma \wedge (\mathbf{n}_s + 3 = \mathbf{n}_r \wedge \mathbf{n}_1 = 2 \vee \mathbf{n}_s + 4 = \mathbf{n}_r \wedge \mathbf{n}_1 = 3))$$

This simplifies the abstraction further. After that, our widening operation compares the current state with the previous one, to look for the same (numerical) constraints that both states imply, and to replace those numerical constraints in the current state with the ones discovered by widening. This operation eventually ensures termination of our analysis. As for the example, some constraints among  $\mathbf{n}_s$ ,  $\mathbf{n}_r$  and  $\mathbf{n}_1$  can be found to make the widened post-state become:

$$(4) \vee (\sigma \wedge \mathbf{n}_s + \mathbf{n}_1 = \mathbf{n}_r - 1 \wedge \mathbf{n}_1 \geq 2) \quad (5)$$

One more iteration of symbolic execution will produce the same result as (5), suggesting that it is already the fixpoint (and hence the loop invariant):

$$\begin{aligned} & \mathbf{r}::\mathbf{sll}\langle \mathbf{n}_r, \mathbf{a}, \mathbf{b} \rangle \wedge \mathbf{srt} = \mathbf{r} \wedge \mathbf{prv} = \mathbf{null} \wedge \mathbf{n}_r + 1 = \mathbf{n} - \mathbf{n}_x \vee \\ & \mathbf{r}::\mathbf{node}\langle \mathbf{a}, \mathbf{srt} \rangle * \mathbf{srt}::\mathbf{sll}\langle \mathbf{n}_s, \mathbf{c}_1, \mathbf{b} \rangle \wedge \mathbf{prv} = \mathbf{r} \wedge \\ & \quad \mathbf{a} \leq \mathbf{c}_1 \wedge \mathbf{a} \leq \mathbf{v} \wedge \mathbf{n}_r + 1 = \mathbf{n} - \mathbf{n}_x \wedge \mathbf{n}_s + 1 = \mathbf{n}_r \vee \\ & \mathbf{r}::\mathbf{node}\langle \mathbf{a}, \mathbf{prv} \rangle * \mathbf{prv}::\mathbf{node}\langle \mathbf{c}_1, \mathbf{srt} \rangle * \mathbf{srt}::\mathbf{sll}\langle \mathbf{n}_s, \mathbf{c}_2, \mathbf{b} \rangle \wedge \\ & \quad \mathbf{a} \leq \mathbf{c}_1 \leq \mathbf{c}_2 \wedge \mathbf{c}_1 \leq \mathbf{v} \wedge \mathbf{n}_r + 1 = \mathbf{n} - \mathbf{n}_x \wedge \mathbf{n}_s + 2 = \mathbf{n}_r \vee \\ & \mathbf{r}::\mathbf{sls}\langle \mathbf{prv}, \mathbf{n}_1, \mathbf{a}, \mathbf{c}_1 \rangle * \mathbf{prv}::\mathbf{node}\langle \mathbf{c}_2, \mathbf{srt} \rangle * \mathbf{srt}::\mathbf{sll}\langle \mathbf{n}_s, \mathbf{c}_3, \mathbf{b} \rangle \wedge \\ & \quad \mathbf{a} \leq \mathbf{c}_1 \leq \mathbf{c}_2 \leq \mathbf{c}_3 \wedge \mathbf{c}_2 \leq \mathbf{v} \wedge \mathbf{n}_r + 1 = \mathbf{n} - \mathbf{n}_x \wedge \mathbf{n}_s + \mathbf{n}_1 = \mathbf{n}_r - 1 \wedge \mathbf{n}_1 \geq 2 \end{aligned}$$

Note that although it is possible to further join the third disjunctive branch with the fourth, our analysis does not do so as it tries to keep the result as precise as possible by eliminating only auxiliary pointer variables.

With the frame part  $\text{cur}::\text{node}\langle v, \mathbf{x} \rangle * \mathbf{x}::\text{ll}\langle \mathbf{n}_x \rangle$  added back, the analysis for the outer loop continues. Eventually, the following loop invariant is discovered for the outer loop:

$$(\mathbf{x}::\text{ll}\langle \mathbf{n}_x \rangle \wedge \mathbf{r}=\text{null} \wedge \mathbf{n}_x=\mathbf{n}) \vee (\mathbf{r}::\text{node}\langle \mathbf{a}, \text{null} \rangle * \mathbf{x}::\text{ll}\langle \mathbf{n}_x \rangle \wedge \mathbf{n}=\mathbf{n}_x+1) \vee (\mathbf{r}::\text{sll}\langle \mathbf{n}_r, \mathbf{a}, \mathbf{b} \rangle * \mathbf{x}::\text{ll}\langle \mathbf{n}_x \rangle \wedge \mathbf{n}=\mathbf{n}_x+\mathbf{n}_r \wedge \mathbf{n}_r \geq 2)$$

which allows us to verify the whole method successfully using e.g. HIP/SLEEK.

### 3 Language and Abstract Domain

To simplify presentation, we focus on a strongly-typed C-like imperative language in Figure 2. The program *Prog* written in this language consists of declarations *tdecl*, which can either be data type declarations *datat* (e.g. `node` in Section 2), or predicate definitions *spred* (e.g. `ll`, `ls`, `sll`, `sls` in Section 2.1), as well as method declarations *meth*. The definitions for *spred* and *mspec* are given later in Figure 3. Without loss of expressiveness, we use an expression-oriented language. So the body of a method (*e*) is an expression formed by standard commands of an imperative language. Note that *d* and *d[v]* represent resp. heap-insensitive and heap sensitive commands. The language allows both call-by-value and call-by-reference method parameters (separated with a semicolon `;`).

<i>Prog</i> ::= <i>tdecl</i> * <i>meth</i> *	<i>tdecl</i> ::= <i>datat</i>   <i>spred</i>
<i>datat</i> ::= <code>data</code> <i>c</i> { <i>field</i> * }	<i>field</i> ::= <i>t v</i> <i>t</i> ::= <i>c</i>   $\tau$
<i>meth</i> ::= <i>t mn</i> (( <i>t v</i> )*; ( <i>t v</i> )*) <i>mspec</i> { <i>e</i> }	$\tau$ ::= <code>int</code>   <code>bool</code>   <code>void</code>
<i>e</i> ::= <i>d</i>   <i>d[v]</i>   <i>v</i> := <i>e</i>   <i>e</i> <sub>1</sub> ; <i>e</i> <sub>2</sub>   <i>t v</i> ; <i>e</i>   <code>if</code> <i>v</i> <code>then</code> <i>e</i> <sub>1</sub> <code>else</code> <i>e</i> <sub>2</sub>   <code>while</code> <i>v</i> { <i>e</i> }	
<i>d</i> ::= <code>null</code>   <i>k</i> <sup><math>\tau</math></sup>   <i>v</i>   <code>new</code> <i>c</i> ( <i>v</i> *)   <i>mn</i> ( <i>u</i> *; <i>v</i> *)	
<i>d[v]</i> ::= <i>v.f</i>   <i>v.f</i> := <i>w</i>   <code>free</code> ( <i>v</i> )	

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

Our specification language (in Figure 3) allows (user-defined) shape predicates *spred* to specify both shape and numerical properties. Note that *spred* are constructed with disjunctive constraints  $\Phi$  and numerical formulae  $\pi$ . We require that the predicates be well-formed [21].

A conjunctive abstract program state,  $\sigma$ , is composed of a heap (shape) part  $\kappa$  and a numerical part  $\pi$ , where  $\pi$  consists of  $\gamma$  and  $\phi$  as aliasing and numerical 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  $\Delta$  (and its set is recognised as  $\mathcal{P}_{\text{SH}}$ ). An abstract state  $\Delta$  can be normalised to the  $\Phi$  form.

Using entailment [21], we define a partial order over these abstract states:

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

$spread$	$::= \text{root}::c(v^*) \equiv \Phi$	$\Phi ::= \bigvee \sigma^*$	$\sigma ::= \exists v^* \cdot \kappa \wedge \pi$
$mspec$	$::= \text{requires } \Phi_{pr} \text{ ensures } \Phi_{po}$		
$\Delta$	$::= \Phi \mid \Delta_1 \vee \Delta_2 \mid \Delta \wedge \pi \mid \Delta_1 * \Delta_2 \mid \exists v \cdot \Delta$		
$\kappa$	$::= \text{emp} \mid v::c(v^*) \mid \kappa_1 * \kappa_2$	$\pi ::= \gamma \wedge \phi$	
$\gamma$	$::= v_1 = v_2 \mid v = \text{null} \mid v_1 \neq v_2 \mid v \neq \text{null} \mid \text{true} \mid \gamma_1 \wedge \gamma_2$		
$\phi$	$::= b \mid a \mid \phi_1 \wedge \phi_2 \mid \phi_1 \vee \phi_2 \mid \neg \phi \mid \exists v \cdot \phi \mid \forall v \cdot \phi$		
$b$	$::= \text{true} \mid \text{false} \mid v \mid b_1 = b_2$	$a ::= s_1 = s_2 \mid s_1 \leq s_2$	
$s$	$::= k^{\text{int}} \mid v \mid k^{\text{int}} \times s \mid s_1 + s_2 \mid -s \mid \max(s_1, s_2) \mid \min(s_1, s_2)$		

**Fig. 3.** The Specification Language.

where  $\mathbf{R}$  is the (computed) residue part. And we also have an induced lattice over these states as the base of fixpoint calculation for loop invariants.

The memory model of our specification formula is similar to the model given for separation logic [25], except that we have extensions to handle user-defined shape predicates and related numerical properties. In our analysis, all the variables except the program ones are logical variables. We denote a program variable's initial value as unprimed and its current value as primed [21].

## 4 Analysis Algorithm

Our proposed analysis algorithm is given in Figure 4.

<b>Fixpoint Computation in Combined Domain</b>	
<b>Input:</b>	$\mathcal{T}, \Delta_{pre}, \text{while } b \{e\}, n;$
<b>Local:</b>	$i := 0; \Delta_i := \text{false}; \Delta'_i := \text{false};$
1	<b>repeat</b>
2	$i := i + 1;$
3	$\Delta_i := \text{widen}^\dagger(\Delta_{i-1}, \text{join}^\dagger(\Delta_{pre}, \Delta'_{i-1}));$
4	$\Delta'_i := \text{abs}^\dagger(\llbracket e \rrbracket_{\mathcal{T}}(\Delta_i \wedge b));$
5	<b>if</b> $\Delta'_i = \text{false} \vee \text{cp\_no}(\Delta'_i) > n$ <b>then return fail end if</b>
6	<b>until</b> $\Delta'_i = \Delta'_{i-1};$
7	<b>return</b> $\Delta'_i$

**Fig. 4.** Main analysis algorithm.

The algorithm takes four input parameters:  $\mathcal{T}$  as the program environment with all the method specifications in the program (for potential method calls in loop body),  $\Delta_{pre}$  as the precondition of the loop's (symbolic) execution, the while loop itself  $\text{while } b \{e\}$ , and the number of upper bound of shared logical variables we keep during the analysis  $n$ .

Our analysis is based on abstract interpretation [7] with specifically designed operations (**abs**, **join** and **widen**) over this combined domain.<sup>2</sup> At the beginning, we initialise the iteration variable ( $i$ ) and two states to begin with ( $\Delta_i$  and  $\Delta'_i$ ). The **false**'s here as initial values denote the top element of our defined lattice

<sup>2</sup> Note that our analysis uses lifted versions of these operations (indicated by  $\dagger$ ), which will be explained in more details in Section 4.2.

as well as our starting point of the fixpoint iteration. Among the two states here, the unprimed version  $\Delta_i$  denotes the initial state before the  $i^{\text{th}}$  execution of the loop body, and the primed one  $\Delta'_i$  represents the result state after. Each iteration starts at line 1. Firstly we join together the precondition of the loop with the result state  $\Delta'_{i-1}$  obtained in the previous iteration, and widen it against the initial state  $\Delta_{i-1}$  of the previous iteration (line 3). Then we symbolically execute the loop body with the abstract semantics in Section 4.1 (line 4), and apply the abstraction operation to the obtained abstract state. If the symbolic execution cannot continue due to a program bug, or if we find our abstraction cannot keep the number of shared logical variables/cutpoints (counted by `cp_no`) within a specified bound ( $n$ ), then a failure is reported (line 5). Otherwise we judge whether a fixpoint is already reached by comparing the current abstract state with the previous one (line 6). The fixpoint  $\Delta'_i$  is returned as the loop invariant.

We will elaborate the key techniques of our analysis in what follows: the abstract semantics, the abstraction function, and the join and widening operators.

#### 4.1 Abstract Semantics

The abstract semantics is used to execute the loop body symbolically to obtain its post-state during the loop invariant synthesis. Its type is defined as

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

where `AllSpec` contains all the specifications of all methods (extracted from the program *Prog*). For some expression  $e$ , given its precondition, the semantics will calculate the postcondition.

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

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

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

$$\frac{\text{isdatat}(c) \quad \sigma \vdash x::c\langle v^* \rangle * \sigma'}{\text{rearr}(x)\sigma =_{df} \sigma} \quad \frac{\text{isspred}(c) \quad \sigma \vdash x::c\langle u^* \rangle * \sigma' \quad \text{root}::c\langle v^* \rangle \equiv \Phi}{\text{rearr}(x)\sigma =_{df} \sigma' * [x/\text{root}, u^*/v^*]\Phi}$$

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.  $x.f_i$ ,  $x.f_i := w$ , or  $\text{free}(x)$ ) assumes that the rearrangement  $\text{rearr}(x)$  has been done in prior:



$$\begin{array}{c}
 \frac{\text{isdatat}(c) \quad \sigma \vdash x::c\langle v_1, \dots, v_n \rangle * \sigma'}{\text{exec}(x.f_i)(\mathcal{T})\sigma =_{df} \sigma' * x::c\langle v_1, \dots, v_n \rangle \wedge \mathbf{res}=v_i} \\
 \frac{\text{isdatat}(c) \quad \sigma \vdash x::c\langle v_1, \dots, v_n \rangle * \sigma'}{\text{exec}(x.f_i := w)(\mathcal{T})\sigma =_{df} \sigma' * x::c\langle v_1, \dots, v_{i-1}, w, v_{i+1}, \dots, v_n \rangle} \\
 \frac{\text{isdatat}(c) \quad \sigma \vdash x::c\langle u^* \rangle * \sigma'}{\text{exec}(\mathbf{free}(x))(\mathcal{T})\sigma =_{df} \sigma'}
 \end{array}$$

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

$$\begin{array}{c}
 \text{exec}(k)(\mathcal{T})\sigma =_{df} \sigma \wedge \mathbf{res}=k \qquad \text{exec}(x)(\mathcal{T})\sigma =_{df} \sigma \wedge \mathbf{res}=x \\
 \frac{\text{isdatat}(c)}{\text{exec}(\mathbf{new } c(v^*))(\mathcal{T})\sigma =_{df} \sigma * \mathbf{res}::c\langle v^* \rangle} \\
 \frac{t \text{ mn } ((t_i \ u_i)_{i=1}^m; (t'_i \ v_i)_{i=1}^n) \text{ requires } \Phi_{pr} \text{ ensures } \Phi_{po} \in \mathcal{T} \quad \rho = [x'_i/u_i]_{i=1}^m \circ [y'_i/v_i]_{i=1}^n \quad \sigma \vdash \rho\Phi_{pr} * \sigma' \quad \rho_o = [r_i/v_i]_{i=1}^n \circ [x'_i/u'_i]_{i=1}^m \circ [y'_i/v'_i]_{i=1}^n \quad \rho_l = [r_i/y'_i]_{i=1}^n \quad \text{fresh logical } r_i}{\text{exec}(\text{mn}(x_1, \dots, x_m; y_1, \dots, y_n))(\mathcal{T})\sigma =_{df} (\rho_l\sigma') * (\rho_o\Phi_{po})}
 \end{array}$$

Note that the first three rules deal with constant ( $k$ ), variable ( $x$ ) and data node creation ( $\mathbf{new } c(v^*)$ ), respectively, while the last rule handles method invocation. 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  $\rho_{ol} \circ \rho_o\Phi_{po}$ .

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

$$\mathbf{rearr}^\dagger(x) \bigvee \sigma_i =_{df} \bigvee (\mathbf{rearr}(x)\sigma_i)$$

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

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

Based on the transition functions above, we can define the abstract semantics for a program command  $e$  as follows:

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

which form the foundation for us to analyse the loop body.

## 4.2 Abstraction, Join and Widening

This section describes our specifically designed abstraction, join and widening operations employed in our loop invariant synthesis process.

**Abstraction function.** During the symbolic execution, we may be confronted with many ‘‘concrete’’ shapes in postconditions of the loop body. As an example

of list traversal, the list may contain one node, or two nodes, or even more nodes in the list, which the analysis cannot enumerate infinitely. The abstraction function deals with those situations by abstracting the (potentially infinite) concrete situations into more abstract shapes. Our rationale is to keep only program variables and shared cutpoints; all other logical variables will be abstracted away. As an instance, the first state below can be further abstracted (as shown), while the second one cannot:

$$\begin{aligned} \text{abs}(\mathbf{x}::\text{node}\langle -, \mathbf{z}_0 \rangle * \mathbf{z}_0::\text{node}\langle -, \text{null} \rangle) &= \mathbf{x}::\text{ll}\langle \mathbf{n} \rangle \wedge \mathbf{n}=2 \\ \text{abs}(\mathbf{x}::\text{node}\langle -, \mathbf{z}_0 \rangle * \mathbf{y}::\text{node}\langle -, \mathbf{z}_0 \rangle * \mathbf{z}_0::\text{node}\langle -, \text{null} \rangle) &= - \end{aligned} \quad (6)$$

where both  $\mathbf{x}$  and  $\mathbf{y}$  are program variables, and  $\mathbf{z}_0$  is an existentially quantified logical variable. In the second case  $\mathbf{z}_0$  is a shared cutpoint referenced by both  $\mathbf{x}$  and  $\mathbf{y}$ , and thus the state is not changed. As illustrated, the abstraction transition function  $\text{abs}$  eliminates unimportant cutpoints (during analysis) to ensure termination. Its type is defined as follows:

$$\text{abs} : \text{SH} \rightarrow \text{SH} \quad \text{Abstraction}$$

which indicates that it takes in a conjunctive abstract state  $\sigma$  and abstracts it as another conjunctive state  $\sigma'$ . Below are its rules.

$$\begin{array}{c} \text{abs}(\sigma \wedge x_0=e) =_{df} \sigma[e/x_0] \qquad \text{abs}(\sigma \wedge e=x_0) =_{df} \sigma[e/x_0] \\ \frac{x_0 \notin \text{Reach}(\sigma)}{\text{abs}(x_0::c\langle v^* \rangle * \sigma) =_{df} \sigma * \text{true}} \\ \frac{\begin{array}{l} \text{isdata}(c_1) \qquad c_2\langle u_2^* \rangle \equiv \Phi \\ p::c_1\langle v_1^* \rangle * \sigma_1 \vdash p::c_2\langle v_2^* \rangle * \sigma_2 \qquad \text{Reach}(p::c_2\langle v_2^* \rangle * \sigma_2) \cap \{v_1^*\} = \emptyset \end{array}}{\text{abs}(p::c_1\langle v_1^* \rangle * \sigma_1) =_{df} p::c_2\langle v_2^* \rangle * \sigma_2} \end{array}$$

The first two rules eliminate logical variables ( $x_0$ ) by replacing them with their equivalent expressions ( $e$ ). The third rule is used to eliminate any garbage (heap part led by a logical variable  $x_0$  unreachable from the other part of the heap) that may exist in the heap. As  $x_0$  is already unreachable from, and not usable by, the program variables, it is safe to treat it as garbage  $\text{true}$ , for example the  $\mathbf{x}_0$  in  $\mathbf{x}::\text{node}\langle -, \text{null} \rangle * \mathbf{x}_0::\text{node}\langle -, \text{null} \rangle$  where only  $\mathbf{x}$  is a program variable.

The last rule of  $\text{abs}$  plays the most significant role which intends to eliminate shape formulae led by logical variables (all variables in  $v_1^*$ ). It tries to fold data nodes up to a predicate node. It confirms that  $c_1$  is a data node definition and  $c_2$  is a predicate. Meanwhile it also ensures that the latter is a sound abstraction of the former by entailment checking, and the logical parameters of  $c_1$  are not reachable from other part of the heap (so that the abstraction does not lose necessary information). The function  $\text{Reach}$  is defined as follows:

$$\text{Reach}(\sigma) =_{df} \bigcup_{v \in \text{fv}(\sigma)} \text{ReachVar}(\kappa \wedge \pi, v) \text{ where } \sigma ::= \exists \mathbf{u}^* \cdot \kappa \wedge \pi$$

returning all variables which are reachable from free variables in the abstract state  $\sigma$ . The function  $\text{ReachVar}(\kappa \wedge \pi, v)$  returns the minimal set of variables satisfying the relationship below:

$$\{v\} \cup \{z_2 \mid \exists z_1, \pi_1 \cdot z_1 \in \text{ReachVar}(\kappa \wedge \pi, v) \wedge \pi = (z_1 = z_2 \wedge \pi_1)\} \cup \{z_2 \mid \exists z_1, \kappa_1 \cdot z_1 \in \text{ReachVar}(\kappa \wedge \pi, v) \wedge \kappa = (z_1 :: c(\dots, z_2, \dots) * \kappa_1)\} \subseteq \text{ReachVar}(\kappa \wedge \pi, v)$$

That is, it is composed of aliases of  $v$  and variables reachable from  $v$ . As in the previous example:  $\text{abs}(x::\text{node}(\_, z_0) * z_0::\text{node}(\_, \text{null})) \rightsquigarrow x::\text{ll}(\mathbf{n}_0) \wedge \mathbf{n}_0=2$ .

During the analysis, we apply the above abstraction rules (following the given order) onto the current abstract state exhaustively until it stabilises. Such convergence is confirmed because the abstract shape domain is finite due to the bounded numbers of variables and predicates, as discussed later.

Finally the lifting function is overloaded for  $\text{abs}$  to lift both its domain and range to disjunctive abstract states  $\mathcal{P}_{\text{SH}}$ :

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

which allows it to be used in the analysis.

**Join operator.** The operator  $\text{join}$  is applied over two conjunctive abstract states, trying to find a common shape as a sound abstraction for both:

$$\begin{aligned} \text{join}(\sigma_1, \sigma_2) =_{df} & \\ & \text{let } \sigma'_1, \sigma'_2 = \text{rename}(\sigma_1, \sigma_2) \text{ in} \\ & \text{match } \sigma'_1, \sigma'_2 \text{ with } (\exists x_1^* \cdot \kappa_1 \wedge \pi_1), (\exists x_2^* \cdot \kappa_2 \wedge \pi_2) \text{ in} \\ & \quad \text{if } \kappa_1 \vdash \kappa_2 * \text{true} \text{ then } \exists x_1^*, x_2^* \cdot \kappa_2 \wedge (\text{join}_\pi(\pi_1, \pi_2)) \\ & \quad \text{else if } \kappa_2 \vdash \kappa_1 * \text{true} \text{ then } \exists x_1^*, x_2^* \cdot \kappa_1 \wedge (\text{join}_\pi(\pi_1, \pi_2)) \\ & \quad \text{else } \sigma_1 \vee \sigma_2 \end{aligned}$$

where the  $\text{rename}$  function prevents naming clashes among logical variables of  $\sigma_1$  and  $\sigma_2$ , by renaming logical variables of same name in the two states with fresh names. For example it will renew  $x_0$ 's name in both states  $\exists x_0 \cdot x_0=0$  and  $\exists x_0 \cdot x_0=1$  to make them  $\exists x_0 \cdot x_0=0$  and  $\exists x_1 \cdot x_1=1$ . After this procedure it judges whether  $\sigma_2$  is an abstraction of  $\sigma_1$ , or the other way round. If either case holds, it regards the shape of the weaker state as the shape of the joined states, and performs joining for numerical formulae with  $\text{join}_\pi(\pi_1, \pi_2)$ , the convex hull operator over numerical domain [23]. Otherwise it keeps a disjunction of the two states (as it would be unsound to join their shapes together in this case). Then we lift this operator for abstract state  $\Delta$  as follows:

$$\text{join}^\dagger(\Delta_1, \Delta_2) =_{df} \text{match } \Delta_1, \Delta_2 \text{ with } (\bigvee_i \sigma_i^1), (\bigvee_j \sigma_j^2) \text{ in } \bigvee_{i,j} \text{join}(\sigma_i^1, \sigma_j^2)$$

which essentially joins all pairs of disjunctions from the two abstract states, and makes a disjunction of them.

**Widening operator.** The finiteness of the shape domain is confirmed by the abstraction function. To ensure the termination of the whole analysis, we still need to guarantee the convergence over the numerical domain. This task is accomplished by the widening operator.

The widening operator  $\text{widen}(\sigma_1, \sigma_2)$  is defined as

$$\begin{aligned} \text{widen}(\sigma_1, \sigma_2) =_{df} & \\ & \text{let } \sigma'_1, \sigma'_2 = \text{rename}(\sigma_1, \sigma_2) \text{ in} \\ & \text{match } \sigma'_1, \sigma'_2 \text{ with } (\exists x_1^* \cdot \kappa_1 \wedge \pi_1), (\exists x_2^* \cdot \kappa_2 \wedge \pi_2) \text{ in} \\ & \quad \text{if } \kappa_1 \vdash \kappa_2 * \text{true} \text{ then } \exists x_1^*, x_2^* \cdot \kappa_2 \wedge (\text{widen}_\pi(\pi_1, \pi_2)) \\ & \quad \text{else } \sigma_1 \vee \sigma_2 \end{aligned}$$

where the `rename` function has the same effect as above. Generally this operator is analogous to `join`; the only difference is that we expect the second operand  $\sigma_2$  is weaker than the first  $\sigma_1$ , such that the widening reflects the trend of such weakening from  $\sigma_1$  to  $\sigma_2$ . In this case it applies the widening operation  $\text{widen}_\pi(\pi_1, \pi_2)$  over the numerical domain [23]. Therefore, based on the widening over conjunctive abstract states, we lift the operator over (disjunctive) abstract states:

$$\text{widen}^\dagger(\Delta_1, \Delta_2) =_{df} \mathbf{match} \Delta_1, \Delta_2 \mathbf{with} (\bigvee_i \sigma_i^1), (\bigvee_j \sigma_j^2) \mathbf{in} \bigvee_{i,j} \text{widen}(\sigma_i^1, \sigma_j^2)$$

which is similar as its counterpart of the `join` operator. These three operations provides termination guarantee while preserving soundness, as the following example demonstrates.

*Example 1 (Abstraction, join and widening).* Assume we have two abstract states,  $\Delta_0 = \mathbf{x}::\mathbf{node}\langle -, \mathbf{x}_0 \rangle * \mathbf{x}_0::\mathbf{node}\langle -, \mathbf{null} \rangle$  and  $\Delta_1 = \mathbf{x}::\mathbf{node}\langle -, \mathbf{x}_0 \rangle * \mathbf{x}_0::\mathbf{node}\langle -, \mathbf{x}_1 \rangle * \mathbf{x}_1::\mathbf{node}\langle -, \mathbf{null} \rangle$ . We would like to discover a sound approximation for both states. Firstly we perform abstractions on both to obtain two abstract states, say,  $\Delta'_0 = \mathbf{x}::\mathbf{ll}\langle \mathbf{n}_0 \rangle \wedge \mathbf{n}_0=2$  and  $\Delta'_1 = \mathbf{x}::\mathbf{ll}\langle \mathbf{n}_0 \rangle \wedge \mathbf{n}_0=3$ . Then these two are joined together according to shape similarity to be  $\Delta''_1 = \mathbf{x}::\mathbf{ll}\langle \mathbf{n}_0 \rangle \wedge (\mathbf{n}_0=2 \vee \mathbf{n}_0=3)$ , which transfers disjunction to numerical domain. Finally the joined state is widened based on the first state  $\Delta'_0$ , yielding a state  $\mathbf{x}::\mathbf{ll}\langle \mathbf{n}_0 \rangle \wedge \mathbf{n}_0 \geq 2$ . It is a sound abstraction of both  $\Delta_0$  and  $\Delta_1$ , and finishes the analysis with one more iteration.

**Soundness and termination.** The soundness of our analysis is ensured by the soundness of the following: the entailment prover [21], the abstract semantics (w.r.t. concrete semantics), the abstraction operation over shapes, and the `join` and widening operators.

**Theorem 1 (Soundness).** *Our analysis is sound due to soundness of entailment checking, abstract semantics, operations of abstraction, join and widening.*

The proof for entailment checking is by structural induction [21]; for abstract semantics is by induction over program constructors; for abstraction follows directly the first two; and for `join` and widening is based on entailment checking and soundness of corresponding numerical operators.

For the termination aspect, we have the result:

**Theorem 2 (Termination).** *The iteration of our fixpoint computation will terminate in finite steps for finite input of program and specification.*

The proof is based on two facts: the finiteness over the shape domain provided by our restriction on cutpoints, and the termination over the numerical domain guaranteed by our widening operator. The first can be proved by claiming the finiteness of all possible abstract states only with the shape information: recalling our analysis algorithm where we set an upper bound  $n$  for shared cutpoints (logical variables) we keep in track of, we know that the program and logical variables preserved in our analysis are finite. Meanwhile all possible shape predicates are limited; therefore all the shape-only abstract states are finite. The second is proved in the abstract interpretation frameworks for numerical domains [23]. These two facts guarantee the convergence of our analysis.

## 5 Experiments and Evaluation

We have implemented a prototype system for evaluation purpose. The prototype system was built in Objective Caml. We used SLEEK [21] as the solver for entailment checking over the heap domain, and Omega constraint solver [24] and Fixcalc solver [23] for join and widening operations in the numerical domain. Our test platform was an Intel Core 2 CPU 2.66GHz system with 8Gb RAM.

Program	Function	Time
<code>create</code>	Creates a list with given length parameter	0.452
<code>ins_sort</code>	Inner loop of Fig. 1	0.824
<code>ins_sort</code>	Outer loop of Fig. 1	4.372
<code>delete</code>	Disposes a list	0.720
<code>traverse</code>	Traverses a list	0.636
<code>append</code>	Appends two lists	0.312
<code>partition</code>	Auxiliary operation used by Quick-sort	1.497
<code>merge</code>	Merges two sorted lists to be one sorted list	1.972
<code>split</code>	Divides a list into two sublists with length difference of at most one	0.354
<code>select</code>	Selects the smallest node of a list	0.692
<code>select_sort</code>	Outer loop of selection sort	4.892
<code>tree_insert</code>	Inserts a node into a binary search tree	1.364
<code>tree_search</code>	Finds a node in a binary search tree	1.294

**Fig. 5.** Selected Experimental Results.

Figure 5 describes the programs with which we performed experiments. The first column denotes the names of the programs. The second column states the programs' functionalities. The last column exhibits the time in second taken by our analysis. As can be seen from their functions, these programs involve recursive data structures such as (sorted) linked lists and binary (search) trees, and employ loops to manipulate these data structures (and some of them even have nested loops). Our target is to verify these programs with the help of our analysis over the loops they invoke, such that user annotations for while loops can be avoided. Our experiments have confirmed that HIP/SLEEK can verify all these programs successfully when supplied with loop invariants discovered by our analysis. According to our experience, these experiments just require the bound of shared cutpoints be a reasonably small number, say no more than twice of the number of program variables.

We have two main observations from our experimental results. The first is that we can handle many data structures with rich program properties they bear. To analyse these loops, we need to deal with both the list and list segment predicates to capture the linked list data structure, as well as their sorted version for the sorting algorithms. We can also handle tree-like predicates such as binary trees and binary search trees. Meanwhile these predicates also come along with many properties such as the length of the list and size/height of the tree, and the

minimum/maximum value of a sorted list/binary search tree. Based on them, our analysis is capable of expressing the invariants of these properties in terms of the constraints over the predicates' parameters.

Beyond the number of predicates and properties we can process, another observation on our analysis is that we can process them *rather precisely*. For example the list creation program creates a list with the same length as user input, and list traverse does not change list's length. Besides these, some loops provide critical invariants for the method running them to function correctly. For example, the quicksort algorithm partitions a list into three parts, where two are lists and the third just one node, whose value is exactly in the middle of that of the two other lists (`partition` in the table). We use a list bound predicate to indicate that fact which is successfully inferred by our analysis. We can also infer that the first loop of a mergesort (`split` in the table) can divide the list into two whose length difference is at most one, which is unimportant for the algorithm's functional correctness but essential for its performance. For `tree_insert`, we have the result that the tree's height is increased at most one, and the minimum/maximum value of the new binary search tree will be exactly the inserted value, if that value is out of the value bounds of the original tree. Such invariants are sufficiently precise to prove the functional correctness of all these programs with the given predicates.

## 6 Related Work and Conclusion

**Related works.** For heap-manipulating programs with any form of recursion (be it loop or recursive method call), dramatic advances have been made in synthesising their invariants/specifications. The local shape analysis [9] infers loop invariants for list-processing programs, followed by the SpaceInvader tool to infer full method specifications over the separation domain, so as to verify pointer safety for larger industrial codes [3, 27]. The SLAyer tool [10] implements an inter-procedural analysis for programs with shape information. To deal with also size information (such as number of nodes in lists/trees), THOR [19] derives a numerical program from the original heap-processing one in a sound way, such that the size information can be obtained with a traditional loop invariant synthesis. A similar approach [11] 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 handle data structures with stronger invariants such as sortedness and binary search property, which have not been addressed in the previous works.

One more work to be mentioned is the relational inductive shape analysis [4]. It employs inductive checkers to express both shape and numerical information. Our approach has four advantages over theirs: firstly, we try to keep as many as possible shared cutpoints (logical variables) during the analysis (within a preset bound), whereas they do not preserve such cutpoints (which is witnessed by their joining rules over the shape domain). Therefore our analysis is essentially more precise than theirs, e.g. in the second scenario of (6) described in Section 4.2.

Meanwhile, our approach can deal with data structures with loops in them (say cyclic linked-lists), whereas they do not have a mechanism to handle it. An example in point is the state  $x::ls\langle m, y \rangle * y::ls\langle y, n \rangle \wedge n > 0$  involving both a shared cutpoint  $y$  and a circled list  $y::ls\langle y, n \rangle \wedge n > 0$ , neither of which can be handled by their work (while ours is capable of that). Another advantage of our approach over theirs is that they only demonstrate how to analyse a program with one particular shape, such as their examples analysing programs which manipulate binary search trees and red-black trees without changing the variety of shape in the heap. Comparatively, we allow different predicates to appear in the analysis of one program, like in our motivating example (thanks to our more flexible abstraction operation). Lastly, their work is mainly from a theory perspective as they do not employ numerical reasoners to solve the relational constraints in their implementation; on the contrary, we discharge all the numerical and relational constraints with automated reasoners [23, 24].

There are also many other approaches that can synthesise shape-related program invariants than those based on separation logic. The shape analysis framework TVLA [26] is based on three-valued logic. It is capable of handling complicated data structures and properties, such as sortedness. Guo et al. [12] reported a global shape analysis that discover inductive structural shape invariants from the code. Kuncak et al. [16] 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. Hackett and Rugina [13] can deal with AVL-trees but is customised to handle only tree-like structures with height property. Compared with these works, separation logic based approach benefits from the frame rule and hence supports local reasoning.

Classical abstract interpretation [7] and its applications such as automated assertion discovery [8, 15, 17] mainly focus on finding numerical program properties. Compared with their works, ours is also founded on the abstract interpretation framework but tries to discover loop invariants with *both* separation and numerical information. Meanwhile, we can also utilise their techniques of join and widening to reason about the numerical domain, as we did for the work [23].

**Concluding Remarks.** We have reported an analysis which allows us to synthesise sound and useful loop invariants over a combined separation and numerical domain. The key components of our analysis include novel operations for abstraction, join and widening in the combined domain. We have built a prototype system and the initial experimental results are encouraging.

**Acknowledgement.** This work was supported by EPSRC Projects EP/G042322/1 and EP/E021948/1 and MoE Tier-2 Project R-252-000-444-112. We thank Florin Craciun for his precious comments.

## References

1. J. Berdine, C. Calcagno, and P. O’Hearn. Smallfoot: Modular automatic assertion checking with separation logic. In *FMCO*, 2005.
2. J. Berdine, B. Cook, D. Distefano, and P. O’Hearn. Automatic termination proofs for programs with shape-shifting heaps. In *18th CAV*, 2006.

3. C. Calcagno, D. Distefano, P. O'Hearn, and H. Yang. Compositional shape analysis by means of bi-abduction. In *36th POPL*, January 2009.
4. B. Chang and X. Rival. Relational inductive shape analysis. In *POPL*, 2008.
5. W.-N. Chin, C. David, H. H. Nguyen, and S. Qin. Automated verification of shape, size and bag properties. In *12th ICECCS*, 2007.
6. W.-N. Chin, C. David, H. H. Nguyen, and S. Qin. Enhancing modular oo verification with separation logic. In *POPL*, January 2008.
7. 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.
8. P. Cousot and R. Cousot. On abstraction in software verification. In *CAV*, 2002.
9. D. Distefano, P. O'Hearn, and H. Yang. A local shape analysis based on separation logic. In *TACAS*, 2006.
10. A. Gotsman, J. Berdine, and B. Cook. Interprocedural shape analysis with separated heap abstractions. In *SAS*, 2006.
11. S. Gulwani, T. Lev-Ami, and M. Sagiv. A Combination Framework for Tracking Partition Sizes. In *POPL*, 2009.
12. B. Guo, N. Vachharajani, and D. August. Shape analysis with inductive recursion synthesis. In *PLDI*, 2007.
13. B. Hackett and R. Rugina. Region-based shape analysis with tracked locations. In *POPL*, 2005.
14. S. Ishtiaq and P. O'Hearn. BI as an assertion language for mutable data structures. In *POPL*, 2001.
15. L. Kovács and T. Jebelean. An algorithm for automated generation of invariants for loops with conditionals. In *SYNASC (Symbolic and Numeric Algorithms for Scientific Computing)*, 2005.
16. V. Kuncak, P. Lam, and M. Rinard. Role analysis. In *POPL*, 2002.
17. K. R. M. Leino and F. Logozzo. Loop invariants on demand. In *APLAS*, 2005.
18. K. R. M. Leino. Dafny: an automatic program verifier for functional correctness. To appear at *LPAR-16*, 2010.
19. S. Magill, M. Tsai, P. Lee, and Y. Tsay. Thor: A tool for reasoning about shape and arithmetic. In *CAV*, 2008.
20. H. H. Nguyen and W.-N. Chin. Enhancing program verification with lemmas. In *20th CAV*, 2008.
21. 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.
22. M. Parkinson and G. Bierman. Separation logic, abstraction and inheritance. In *POPL*, 2008.
23. C. Popeea and W.-N. Chin. Inferring disjunctive postconditions. In *Proceedings of 11th Asian Computing Science Conference*, 2006.
24. P. Pugh. The Omega Test: a fast and practical integer programming algorithm for dependence analysis. In *Communications of the ACM*, 1992.
25. J. Reynolds. Separation logic: a logic for shared mutable data structures. In *17th LICS*, 2002.
26. 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.
27. 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.