Automatic Verification of Pointer Programs Using Grammar-Based Shape Analysis^{*}

Oukseh Lee¹, Hongseok Yang², and Kwangkeun Yi³

¹ Dept. of Computer Science & Engineering, Hanyang University, Korea
² ERC-ACI, Seoul National University, Korea

 $^3\,$ School of Computer Science & Engineering, Seoul National University, Korea

Abstract. We present a program analysis that can automatically discover the shape of complex pointer data structures. The discovered invariants are, then, used to verify the absence of safety errors in the program, or to check whether the program preserves the data consistency. Our analysis extends the shape analysis of Sagiv et al. with grammar annotations, which can precisely express the shape of complex data structures. We demonstrate the usefulness of our analysis with binomial heap construction and the Schorr-Waite tree traversal. For a binomial heap construction algorithm, our analysis returns a grammar that precisely describes the shape of a binomial heap; for the Schorr-Waite tree traversal, our analysis shows that at the end of the execution, the result is a tree and there are no memory leaks.

1 Introduction

We show that a static program analysis can automatically verify pointer programs, such as binomial heap construction and the Schorr-Waite tree traversal. The verified properties are: for a binomial heap construction algorithm, our analysis verifies that the returned heap structure is a binomial heap; for the Schorr-Waite tree traversal, it verifies that the output tree is a binary tree, and there are no memory leaks. In both cases, the analysis took less than 0.2 second in Intel Pentium 3.0C with 1GB memory, and its result is simple and human-readable.

Note that although these programs handle regular heap structures such as binomial heaps and trees, the topology of pointers (e.g., cycles) and their imperative operations (e.g., pointer swapping) are fairly challenging for fully automatic static verification without any annotation from the programmer.

The static analysis is an extension to Sagiv et al.'s shape analysis [13] by grammars. To improve accuracy, we associate grammars, which finitely summarize run-time heap structures, with the summary nodes of the shape graphs. This enrichment of shape graph by grammars provides an ample space for precisely capturing the imperative effects on heap structures. The grammar is unfolded to expose an exact heap structure on demand. The grammar is also folded to replace an exact heap structure by an abstract nonterminal. To ensure the termination of the analysis, the grammar merges multiple

^{*} Lee and Yi were supported by the Brain Korea 21 project in 2004, and Yang was supported by R08-2003-000-10370-0 from the Basic Research Program of the Korea Science & Engineering Foundation.

production rules into a single one, and unify multiple nonterminals; this simplification makes the grammar size remain within a practical bound.

The analysis's correctness is proved via separation logic [12, 11]. The analysis is a composition of abstract operations over the grammar-based shape graphs. The semantics (concretization) of the shape graphs is defined as assertions in separation logic. Each abstract operator is proved safe by showing that the separation-logic assertion for the input graph implies that for the output graph. The input program C wrapped by the input and output assertions $\{P\}C\{Q\}$ is always a provable Hoare triple by the separation-logic proof rules.

The main limitation of our analysis is that the analysis cannot handle DAGs and general graphs. To overcome this limitation, we need to use a more general grammar, where the nonterminals can talk about shared cells.

Related Work We borrowed several interesting ideas from the shape analysis [14]. Our analysis represents a program invariant using a set of shape graphs where each shape graph consists of either concrete or abstract nodes. It uses the idea of refining an abstract node, often called focus or materialization, and also the idea of merging the shape graphs which have a similar structure [9, 5].

The difference is the use of grammar; it is the main reason for the improved precision of our analysis. Another difference is that our analysis separates node-summarizing criteria from the properties of the summary nodes. Normally, the shape analysis of Sagiv et al. partitions all the concrete nodes according to the instrumentation predicates that they satisfy, and groups each partition into a single summary node. Thus, two different summary nodes must satisfy different sets of instrumentation predicates. Our analysis, on the other hand, groups the concrete nodes using the most approximate grammar: each group is a maximal set of concrete nodes that can be expressed by the most approximate grammar. Then, the analysis summarizes each group by a single summary node, and annotates the summary node with a new grammar that "best" describes the pointer structure of the summarized concrete nodes. As a consequence, two different summary nodes in our analysis can have the identical grammar annotations.

Graph type [7, 10] and shape type [6] are also closely related to our work. Both of them express invariants of heap objects (or data structures) using grammar-based languages, which are more expressive than the grammars we used. However, they assume that all the loop invariants of a program are provided, while our work infers such invariants.

Outline Section 2 describes the source programming language. Section 3 overviews separation logic that we use to give the meaning of abstract values. Then, we explain the key ideas of our analysis, using a simpler version that can handle tree-like structures with no shared nodes. Section 4 and 5 explain the abstract domain and abstract operators, and Section 6 defines the analyzer. The simpler version is extended to the full analysis in Section 7. Section 8 demonstrates the accuracy of our analysis using binomial heap construction algorithm and the Schorr-Waite tree traversing algorithm.

2 Programming Language

We use the standard while language with additional pointer operations.

This language assumes that every heap cell is binary, having fields 0 and 1. A heap cell is allocated by x := new, and the contents of such an allocated cell is accessed by the field-dereference operation ->. All the other constructs in the language are standard.

3 Separation Logic with Recursive Predicates

Let *Loc* and *Val* be unspecified infinite sets such that nil \notin *Loc* and *Loc* \cup {nil} \subseteq *Val*. We consider separation logic for the following semantic domains.

 $Stack \stackrel{\Delta}{=} Vars \rightharpoonup_{fin} Val$ $Heap \stackrel{\Delta}{=} Loc \rightharpoonup_{fin} Val \times Val$ $State \stackrel{\Delta}{=} Stack \times Heap$

This domain implies that a state has the stack and heap components, and that the heap component of a state has finitely many binary cells.

The assertion language in separation logic is given by the following grammar:⁴

$$P ::= E = E \mid \mathsf{emp} \mid (E \mapsto E, E) \mid P * P \mid \mathsf{true} \mid P \land P \mid P \lor P \mid \neg P \mid \forall x. P$$

Separating conjunction P * Q is the most important, and it expresses the splitting of heap storage; P * Q means that the heap can be split into two parts, so that Pholds for the one and Q for the other. We often use precise equality and iterated separating conjunction, both of which we define as syntactic sugars. Let X be a finite set $\{x_1, \ldots, x_n\}$ where all x_i 's are different.

$$E \doteq E' \stackrel{\Delta}{=} E = E' \land \operatorname{emp} \qquad \bigodot_{x \in X} A_x \stackrel{\Delta}{=} \text{ if } (X = \emptyset) \text{ then emp else } (A_{x_1} * \ldots * A_{x_n})$$

In this paper, we use the extension of the basic assertion language with recursive predicates [15]:

$$P ::= \dots \mid \alpha(E, \dots, E) \mid \text{rec } \Gamma \text{ in } P \qquad \Gamma ::= \alpha(x_1, \dots, x_n) = P \mid \Gamma, \Gamma$$

The extension allows the definition of new recursive predicates by least-fixed points in "rec Γ in P", and the use of such defined recursive predicates in $\alpha(E, \ldots, E)$. To ensure the existence of the least-fixed point in rec Γ in P, we will consider only well-formed Γ where all recursively defined predicates appear in positive positions.

A recursive predicate in this extended language means a set of heap objects. A *heap object* is a pair of location (or locations) and heap. Intuitively, the first component denotes the starting address of a data structure, and the second the cells in the data structure. For instance, when a linked list is seen as a heap object, the location of the head of the list becomes the first component, and the cells in the linked list the second component.

The precise semantics of this assertion language is given by a forcing relation \models . For a state (s,h) and an environment η for recursively defined predicates, we define inductively when an assertion P holds for (s,h) and η . We show the sample clauses below; the full definition appears in [8].

⁴ The assertion language also has the adjoint -* of *. But this adjoint is not used in this paper, so we omit it here.

$$\begin{array}{ll} (s,h),\eta \models \alpha(E) & \text{iff} \quad (\llbracket E \rrbracket s,h) \in \eta(\alpha) \\ (s,h),\eta \models \mathsf{rec} \; \alpha(x) = P \; \text{in} \; Q \quad \text{iff} \quad (s,h),\eta[\alpha \to k] \models P \\ (\text{where} \; k = \mathrm{lfix} \; \lambda k_0.\{(v',h') \mid (s[x \to v'],h'),\eta[\alpha \to k_0] \models P\}) \end{array}$$

4 Abstract Domain

Shape Graph Our analysis interprets a program as a (nondeterministic) transformer of *shape graphs*. A shape graph is an abstraction of concrete states; this abstraction maintains the basic "structure" of the state, but abstracts away all the other details. For instance, consider a state ($[x-1, y-3], [1 \rightarrow \langle 2, \operatorname{nil} \rangle, 2 \rightarrow \langle \operatorname{nil}, \operatorname{nil} \rangle, 3 \rightarrow \langle 1, 3 \rangle$]). We obtain a shape graph from this state in two steps. First, we replace the specific addresses, such as 1 and 2, by *symbolic locations*; we introduce symbols a, b, c, and represent the state by ($[x \rightarrow a, y \rightarrow c], [a \rightarrow \langle b, \operatorname{nil} \rangle, b \rightarrow \langle \operatorname{nil}, \operatorname{nil} \rangle, c \rightarrow \langle a, c \rangle$]). Note that this process abstracts away the specific addresses and just keeps the relationship between the addresses. Second, we abstract heap cells a and b by a grammar. Thus, this step transforms the state to ($[x \rightarrow a, y \rightarrow c], [a \rightarrow \operatorname{tree}, c \rightarrow \langle a, c \rangle$]) where $a \rightarrow \operatorname{tree}$ means that a is the address of the root of a tree, whose structure is summarized by grammar rules for nonterminal tree.

The formal definition of a shape graph is given as follows:

$$\begin{aligned} SymLoc &\stackrel{\Delta}{=} \{a, b, c, \ldots\} \qquad NonTerm \stackrel{\Delta}{=} \{\alpha, \beta, \gamma, \ldots\} \\ Graph \stackrel{\Delta}{=} (Vars \rightharpoonup_{fin} SymLoc) \times (SymLoc \rightharpoonup_{fin} \{nil\} + SymLoc^2 + NonTerm) \end{aligned}$$

Here the set of nonterminals is disjoint from *Vars* and *SymLoc*; these nonterminals represent recursive heap structures such as tree or list. Each shape graph has two components (s, g). The first component s maps stack variables to symbolic locations. The other component g describes heap cells reachable from each symbolic location. For each a, either no heap cells can be reached from a, i.e, g(a) = nil; or, a is a binary cell with contents $\langle b, c \rangle$; or, the cells reachable from a form a heap object specified by a nonterminal α . We also require that g describes all the cells in the heap; for instance, if g is the empty partial function, it means the empty heap.

The semantics (or concretization) of a shape graph (s, g) is given by a translation into an assertion in separation logic:

$$\begin{split} & \mathsf{means}_\mathsf{v}(a,\mathsf{nil}) \mathop{\stackrel{\Delta}{=}} a \doteq \mathsf{nil} \quad \mathsf{means}_\mathsf{v}(a,\alpha) \mathop{\stackrel{\Delta}{=}} \alpha(a) \quad \mathsf{means}_\mathsf{v}(a,\langle b,c\rangle) \mathop{\stackrel{\Delta}{=}} (a \mapsto b,c) \\ & \mathsf{means}_\mathsf{s}(s,g) \mathop{\stackrel{\Delta}{=}} \exists a. \; (\bigodot_{x \in \mathsf{dom}(s)} x \doteq s(x)) * (\bigodot_{a \in \mathsf{dom}(g)} \mathsf{means}_\mathsf{v}(a,g(a))) \end{split}$$

The translation function means_s calls a subroutine means_v to get the translation of the value of g(a), and then, it existentially quantifies all the symbolic locations appearing in the translation. For instance, $\mathsf{means}_s([x \to a, y \to c], [a \to \mathsf{tree}, c \to \langle a, c \rangle])$ is $\exists ac. (x \doteq a) * (y \doteq c) * \mathsf{tree}(a) * (c \mapsto a, c).$

When we present a shape graph, we interchangeably use the set notation and a graph picture. Each variable or symbolic location becomes a node in a graph, and s and g are represented by edges or annotations. For instance, we draw a shape graph $(s,g) = ([x \rightarrow a], [a \rightarrow \langle b, c \rangle, b \rightarrow \operatorname{nil}, c \rightarrow \alpha])$ as:



Note that pair g(a) is represented by two edges (the left one is for field 0 and the right one for field 1), and non-pair values g(b) and g(c) by annotations to the nodes.

Grammar A grammar gives the meaning of nonterminals in a shape graph. We define a grammar R as a finite partial function from nonterminals (the lhs of production rules) to $\wp_{nf}({nil} + ({nil} + NonTerm)^2)$ (the rhs of production rules), where $\wp_{nf}(X)$ is the family of all nonempty finite subsets of X.

$$Grammar \stackrel{\Delta}{=} NonTerm \rightharpoonup_{fin} \wp_{nf}(\{nil\} + (\{nil\} + NonTerm)^2)$$

Set $R(\alpha)$ contains all the possible shapes of heap objects for α . If nil $\in R(\alpha)$, α can be the empty heap object. If $\langle \beta, \gamma \rangle \in R(\alpha)$, then some heap object for α can be split into a root cell, the left heap object β , and the right heap object γ . For instance, if R(tree) = $\{\text{nil}, \langle \text{tree}, \text{tree} \rangle\}$ (i.e., in the production rule notation, tree ::= nil | $\langle \text{tree}, \text{tree} \rangle$), then tree represents binary trees. In our analysis, we use only *well-formed* grammars, where all nonterminals appearing in the range of a grammar are defined in the grammar.

The meaning $\operatorname{means}_{g}(R)$ of a grammar R is given by a recursive predicate declaration Γ in separation logic. Γ is defined exactly for dom(R), and satisfies the following: when nil $\notin R(\alpha)$, $\Gamma(\alpha)$ is

$$\alpha(a) = \bigvee_{\langle v,w\rangle \in R(\alpha)} \exists bc. \ (a \mapsto b,c) * \mathsf{means}_\mathsf{v}(b,v) * \mathsf{means}_\mathsf{v}(c,w),$$

where neither b nor c appears in a, v or w; otherwise, $\Gamma(\alpha)$ is identical as above except that $a \doteq \operatorname{nil}$ is added as a disjunct. For instance, $\operatorname{means}_{g}([\operatorname{tree}_{\operatorname{nil}}, \langle \operatorname{tree}, \operatorname{tree} \rangle])$ is $\{\operatorname{tree}(a) = a \doteq \operatorname{nil} \lor \exists bc. (a \mapsto b, c) * \operatorname{tree}(b) * \operatorname{tree}(c) \}.$

Abstract Domain The abstract domain \widehat{D} for our analysis consists of pairs of a shape graph set and a grammar: $\widehat{D} \triangleq \{\top\} + \wp_{nf}(Graph) \times Grammar$. The element \top indicates that our analysis fails to produce any meaningful results for a given program because the program has safety errors, or the program uses data structures too complex for our analysis to capture. The meaning of each abstract state (\mathcal{G}, R) in \widehat{D} is means $(\mathcal{G}, R) \triangleq$ rec means_g(R) in $\bigvee_{(s,g)\in\mathcal{G}}$ means_s(s,g).

5 Normalized Abstract States and Normalization

The main strength of our analysis is to automatically discover a grammar that describes, in an "intuitive" level, invariants for heap data structures, and to abstract concrete states according to this discovered grammar. This inference of high-level grammars is mainly done by the normalization function from \hat{D} to a subdomain \hat{D}^{∇} of normalized abstract states. In this section, we explain these two notions, normalized abstract states and normalization function.

5.1 Normalized Abstract States

An abstract state (\mathcal{G}, R) is normalized if it satisfies the following two conditions. First, all the shape graphs (s, g) in \mathcal{G} are abstract enough: all the recognizable heap objects are replaced by nonterminals. Note that this condition on (\mathcal{G}, R) is about individual shape graphs in \mathcal{G} . We call a shape graph *normalized* if it satisfies this condition. Second, an abstract state does not have redundancies: all shape graphs are not *similar*, and all nonterminals have *non-similar* definitions.



Fig. 1. Examples of the Normalized Shape Graphs and Similar Shape Graphs.

Normalized Shape Graphs A shape graph is normalized when it is "maximally" folded. A symbolic location a is *foldable* in (s, g) if g(a) is a pair and there is no path from a to a *shared* symbolic location that is referred more than once. When dom(g) of a shape graph (s, g) does not have any foldable locations, we say that (s, g) is *normalized*. For instance, Figure 1.(a) is not normalized, because b is foldable: b is a pair and does not reach any shared symbolic locations. On the other hand, Figure 1.(b) is normalized, because all the pairs in the graph (i.e., a and c) can reach shared symbolic location e.

Similarity We define three notions of similarity: one for shape graphs, another for two cases of the grammar definitions, and the third for the grammar definitions of two nonterminals.

Two shape graphs are similar when they have the similar structures. Let S be a substitution that renames symbolic locations. Two shape graphs (s, g) and (s', g') are similar up to S, denoted $(s, g) \sim_S^G (s', g')$, if and only if

- 1. $\operatorname{dom}(s) = \operatorname{dom}(s')$ and $S(\operatorname{dom}(g)) = \operatorname{dom}(g')$;
- 2. for all $x \in \text{dom}(s)$, S(s(x)) = s'(x); and
- 3. for all $a \in \text{dom}(g)$, if g(a) is a pair $\langle b, c \rangle$ for some b and c, then $g'(S(a)) = \langle S(b), S(c) \rangle$; if g(a) is not a pair, neither is g'(S(a)).

Intuitively, two shape graphs are S-similar, when equating nil and all nonterminals makes the graphs identical up to renaming S. We say that (s,g) and (s',g') are similar, denoted $(s,g) \sim (s',g')$, if and only if there is a renaming relation S such that $(s,g) \sim_S^G (s',g')$. For instance, in Figure 1.(c), (s_1,g_1) and (s_2,g_2) are not similar because we cannot find a renaming substitution S such that $S(s_1(x)) = S(s_1(y))$ (condition 2). However, (s_1,g_1) and (s_3,g_3) are similar because a renaming substitution $\{d/a, e/b, f/c\}$ makes (s_1,g_1) identical to (s_3,g_3) when nil and all nonterminals are erased from the graphs.

Cases e_1 and e_2 in the grammar definitions are similar, denoted $e_1 \sim^C e_2$, if and only if either both e_1 and e_2 are pairs, or they are both non-pair values. The similarity $E_1 \sim^D E_2$ between grammar definitions E_1 and E_2 uses this case similarity: $E_1 \sim^D E_2$ if and only if, for all cases e in E_1 , E_2 has a similar case e' to e ($e \sim^C e'$), and vice versa. For example, in the grammar

 $\alpha ::= \langle \beta, \operatorname{nil} \rangle, \quad \beta ::= \operatorname{nil} \mid \langle \beta, \operatorname{nil} \rangle, \quad \gamma ::= \langle \gamma, \gamma \rangle \mid \langle \alpha, \operatorname{nil} \rangle$

the definitions of α and γ are similar because both $\langle \gamma, \gamma \rangle$ and $\langle \alpha, \text{nil} \rangle$ are similar to $\langle \beta, \text{nil} \rangle$. But the definitions of α and β are not similar since α does not have a case similar to nil.

Definition 1 (Normalized Abstract States). An abstract state (\mathcal{G}, R) is normalized if and only if

- 1. all shape graphs in \mathcal{G} are normalized;
- 2. for all (s_1, g_1) , $(s_2, g_2) \in \mathcal{G}$, we have $(s_1, g_1) \sim (s_2, g_2) \Rightarrow (s_1, g_1) = (s_2, g_2)$; 3. for all $\alpha \in \operatorname{dom}(R)$ and all cases $e_1, e_2 \in R(\alpha)$, $e_1 \sim^C e_2$ implies that $e_1 = e_2$; 4. for all α, β in $\operatorname{dom}(R)$, $R(\alpha) \sim^D R(\beta)$ implies that $\alpha = \beta$.

We write \widehat{D}^{∇} for the set of normalized abstract states.

k-Bounded Normalized States Unfortunately, the normalized abstract domain \widehat{D}^{∇} does not ensure the termination of the analysis, because it has infinite chains. For each number k, we say that an abstract state (\mathcal{G}, R) is k-bounded iff all the shape graphs in \mathcal{G} use at most k symbolic locations, and we define \widehat{D}_k^{∇} to be the set of k-bounded normalized abstract states. This finite domain \widehat{D}_{k}^{∇} is used in our analysis.

Normalization Function 5.2

The normalize function transforms (\mathcal{G}, R) to a normalized (\mathcal{G}', R') with a further abstraction (i.e., $\mathsf{means}(\mathcal{G}, R) \Rightarrow \mathsf{means}(\mathcal{G}', R')$).⁵ It is defined by the composition of five subroutines: normalize = bound^k \circ simplify \circ unify \circ fold \circ rmjunk.

The first subroutine rmjunk removes all the "imaginary" sharing and garbage due to constant symbolic locations, so that it makes the real sharing and garbage easily detectable in syntax. The subroutine rmjunk applies the following two rules until an abstract state does not change. In the definition, "⊎" is a disjoint union of sets, and "." is a union of partial maps with disjoint domains.

(alias) $(\mathcal{G} \uplus \{(s \cdot [x \rightarrow a], g \cdot [a \rightarrow \operatorname{nil}])\}, R) \rightsquigarrow (\mathcal{G} \cup \{(s \cdot [x \rightarrow a'], g \cdot [a \rightarrow \operatorname{nil}, a' \rightarrow \operatorname{nil}])\}, R)$ where a should appear in (s, g) and a' is fresh.

 $(\mathcal{G} \uplus \{(s, q \cdot [a \rightarrow \text{nil}])\}, R) \rightsquigarrow (\mathcal{G} \cup \{(s, q)\}, R)$ where a does not appear in (s, q) (\mathbf{gc})

For instance, given a shape graph $([x \rightarrow a, y \rightarrow a], [a \rightarrow \text{nil}, c \rightarrow \text{nil}]), (\mathbf{gc})$ collects the "garbage" c, and (alias) eliminates the "imaginary sharing" between x and y by renaming a in $y \rightarrow a$. So, the shape graph becomes $([x \rightarrow a, y \rightarrow b], [a \rightarrow nil, b \rightarrow nil])$.

The second subroutine fold converts a shape graph to a normal form, by replacing all foldable symbolic locations by nonterminals. The subroutine fold repeatedly applies the following rule until the abstract state does not change:

(fold) $(\mathcal{G} \uplus \{(s, g \cdot [a \to \langle b, c \rangle, b \to v, c \to v'])\}, R) \rightsquigarrow (\mathcal{G} \cup \{(s, g \cdot [a \to \alpha])\}, R \cdot [\alpha \to \{\langle v, v' \rangle\}])$ where neither b nor c appears in (s, g), α is fresh, and v and v' are not pairs.

The rule recognizes that the symbolic locations b and c are accessed only via a. Then, it represents cell a, plus the reachable cells from b and c by a nonterminal α . Figure 2 shows how the (fold) folds a tree.

The third subroutine unify merges two similar shape graphs in \mathcal{G} . Let (s, g) and (s',g') be similar shape graphs by the identity renaming Δ (i.e., $(s,g) \sim_{\Delta}^{G} (s',g')$). Then, these two shape graphs are almost identical; the only exception is when g(a)and g'(a) are nonterminals or nil. unify eliminates all such differences in two shape graphs; if g(a) and g'(a) are nonterminals, then unify changes g and g', so that they map a to the same fresh nonterminal γ , and then it defines γ to cover both α and β . The unify procedure applies the following rules to an abstract state (\mathcal{G}, R) until the abstract state does not change:

⁵ The normalize function is a reminiscent of the widening in [2, 3].



Fig. 2. Examples of (fold), (unify), and (unil).

- $\begin{aligned} & (\mathbf{unify}) \; (\mathcal{G} \uplus \{ (s_1, g_1 \cdot [a_1 \rightarrow \alpha_1]), (s_2, g_2 \cdot [a_2 \rightarrow \alpha_2]) \}, \; R) \\ & \sim (\mathcal{G} \cup \{ (S(s_1), S(g_1) \cdot [a_2 \rightarrow \beta]), (s_2, g_2 \cdot [a_2 \rightarrow \beta]) \}, \; R \cdot [\beta \rightarrow R(\alpha_1) \cup R(\alpha_2)]) \\ & \text{where} \; (s_1, g_1 \cdot [a_1 \rightarrow \alpha_1]) \sim^G_S (s_2, g_2 \cdot [a_2 \rightarrow \alpha_2]), \; S(a_1) \equiv a_2, \; \alpha_1 \not\equiv \alpha_2, \; \text{and} \; \beta \; \text{is fresh.} \end{aligned}$
- (unil) $(\mathcal{G} \uplus \{(s_1, g_1 \cdot [a_1 \to \alpha]), (s_2, g_2 \cdot [a_2 \to \operatorname{nil}])\}, R)$ $\rightsquigarrow (\mathcal{G} \cup \{(S(s_1), S(g_1) \cdot [a_2 \to \beta]), (s_2, g_2 \cdot [a_2 \to \beta])\}, R \cdot [\beta \to R(\alpha) \cup \{\operatorname{nil}\}])$ where $(s_1, g_1 \cdot [a_1 \to \alpha]) \sim_S^G (s_2, g_2 \cdot [a_2 \to \operatorname{nil}]), S(a_1) \equiv a_2$, and β is fresh.

The (**unify**) rule recognizes two similar shape graphs that have different nonterminals at the same position, and replaces those nonterminals by fresh nonterminal β that covers the two nonterminals. The (**unil**) rule deals with the two similar graphs that have, respectively, nonterminal and nil at the same position. For instance, in Figure 2.(b), the left two shape graphs are unified by (**unify**) and (**unil**). We first replace the left children α and β by γ that covers both; that is, to a given grammar R, we add $[\gamma \rightarrow R(\alpha) \cup R(\beta)]$. Then we replace the right children β and nil by δ that covers both.

The fourth subroutine simplify reduces the complexity of grammar by combining similar cases or similar definitions.⁶ It applies three rules repeatedly:

- If the definition of a nonterminal has two similar cases $\langle \beta, v \rangle$ and $\langle \beta', v' \rangle$, and β and β' are different nonterminals, unify nonterminals β and β' . Apply the same for the second field.
- If the definition of a nonterminal has two similar cases $\langle \beta, v \rangle$ and $\langle \operatorname{nil}, v' \rangle$, add the nil case to $R(\beta)$ and replace $\langle \operatorname{nil}, v' \rangle$ by $\langle \beta, v' \rangle$. Apply the same for the second field.
- If the definitions of two nonterminals are similar, unify the nonterminals.

Formally, the three rules are:

- (case) $(\mathcal{G}, R) \rightsquigarrow (\mathcal{G}, R) \{\beta/\alpha\}$ where $\{\langle \alpha, v \rangle, \langle \beta, v' \rangle\} \subseteq R(\gamma)$ and $\alpha \not\equiv \beta$. (same for the second field)
- $\begin{array}{ll} (\mathbf{nil}) & (\mathcal{G}, R \cdot [\alpha \rightarrow E \uplus \{ \langle \beta, v \rangle, \langle \mathrm{nil}, v' \rangle \}]) \rightsquigarrow (\mathcal{G}, R' [\beta \rightarrow R'(\beta) \cup \{\mathrm{nil}\}]) \\ & \text{where } R' = R \cdot [\alpha \rightarrow E \cup \{ \langle \beta, v \rangle, \langle \beta, v' \rangle \}]. \text{ (same for the second field)} \\ (\mathbf{def}) & (\mathcal{G}, R) \rightsquigarrow (\mathcal{G}, R) \{ \beta / \alpha \} \text{ where } R(\alpha) \sim R(\beta) \text{ and } \alpha \not\equiv \beta. \end{array}$

Here, $(\mathcal{G}, R)\{\alpha/\beta\}$ substitutes α for β , and in addition, it removes the definition of β from R and re-defines α such that α covers both α and β :

⁶ The simplify subroutine is similar to the widening operator in [4].

$$(\mathcal{G}, R \cdot [\alpha \to E_1, \beta \to E_2]) \{\alpha/\beta\} \stackrel{\Delta}{=} (\mathcal{G} \{\alpha/\beta\}, R \{\alpha/\beta\} \cdot [\alpha \to (E_1 \cup E_2) \{\alpha/\beta\}]).$$

For example, consider the following transitions:

$$\begin{array}{l} \alpha::=\operatorname{nil} \mid \langle \beta, \beta \rangle \mid \langle \gamma, \gamma \rangle, \beta::= \langle \gamma, \gamma \rangle, \gamma::= \langle \operatorname{nil}, \operatorname{nil} \rangle \\ \stackrel{(\operatorname{case})}{\rightarrow} \alpha::=\operatorname{nil} \mid \langle \beta, \beta \rangle, \beta::= \langle \beta, \beta \rangle \mid \langle \operatorname{nil}, \operatorname{nil} \rangle \stackrel{(\operatorname{nil})}{\xrightarrow{}} \alpha::=\operatorname{nil} \mid \langle \beta, \beta \rangle, \beta::= \langle \beta, \beta \rangle \mid \langle \beta, \operatorname{nil} \rangle \mid \operatorname{nil} \\ \stackrel{(\operatorname{nil})}{\xrightarrow{}} \alpha::=\operatorname{nil} \mid \langle \beta, \beta \rangle, \beta::= \langle \beta, \beta \rangle \mid \operatorname{nil} \stackrel{(\operatorname{def})}{\xrightarrow{}} \alpha::=\operatorname{nil} \mid \langle \alpha, \alpha \rangle \end{array}$$

In the initial grammar, α 's definition has the similar cases $\langle \beta, \beta \rangle$ and $\langle \gamma, \gamma \rangle$, so we apply $\{\beta/\gamma\}$ (**case**). In the second grammar, β 's definition has similar cases $\langle \beta, \beta \rangle$ and $\langle \operatorname{nil}, \operatorname{nil} \rangle$. Thus, we replace nil by β , and add the nil case to β 's definition (**nil**). We apply (**nil**) once more for the second field. In the fourth grammar, since α and β have similar definitions, we apply $\{\alpha/\beta\}$ (**def**). As a result, we obtain the last grammar which says that α describes binary trees.

The last subroutine **bound**^k checks the number of symbolic locations in each shape graph. The subroutine **bound**^k simply gives \top when one of shape graphs has more than k symbolic locations, thereby ensuring the termination of the analysis.⁷

bound^k(\mathcal{G}, R) = if (no (s, g) in \mathcal{G} has more than k symbolic locations) then (\mathcal{G}, R) else \top

Lemma 1. Given every abstract state (\mathcal{G}, R) , normalize (\mathcal{G}, R) always terminates, and its result is a k-bounded normalized abstract state.

6 Analysis

Our analyzer (defined in Figure 3) consists of two parts: the "forward analysis" of commands C, and the "backward analysis" of boolean expressions B. Both of these interpret C and B as functions on abstract states, and they accomplish the usual goals in the static analysis: for an initial abstract state (\mathcal{G}, R) , $[\![C]\!](\mathcal{G}, R)$ approximates the possible output states, and $[\![B]\!](\mathcal{G}, R)$ denotes the result of pruning some states in (\mathcal{G}, R) that do not satisfy B.

One particular feature of our analysis is that the analysis also checks the absence of memory errors, such as null-pointer dereference errors. Given a command C and an abstraction (\mathcal{G}, R) for input states, the result $\llbracket C \rrbracket (\mathcal{G}, R)$ of analyzing the command Ccan be either some abstract state (\mathcal{G}', R') or \top . (\mathcal{G}', R') means that all the results of C from (\mathcal{G}, R) are approximated by (\mathcal{G}', R') , but in addition to this, it also means that no computations of C from (\mathcal{G}, R) can generate memory errors. \top , on the other hand, expresses the possibility of memory errors, or indicates that a program uses the data structures whose complexity goes beyond the current capability of the analysis.

The analyzer unfolds the grammar definition by calling the subroutine unfold. Given a shape graph (s, g), a variable x and a grammar R, the subroutine unfold first checks whether g(s(x)) is a nonterminal or not. If g(s(x)) is a nonterminal α , unfold looks up the definition of α in R and unrolls this definition in the shape graph (s, g): for each case e in $R(\alpha)$, it updates g by $[s(x) \rightarrow e]$. For instance, when $R(\beta) = \{\langle \beta, \gamma \rangle, \langle \delta, \delta \rangle\}$, unfold($([x \rightarrow a], [a \rightarrow \beta]), R, x)$ is shape-graph set $\{([x \rightarrow a], [a \rightarrow \langle \beta, \gamma \rangle]), ([x \rightarrow a], [a \rightarrow \langle \delta, \delta \rangle])\}$.

⁷ Limiting the number of symbolic locations to be at most k ensures the termination of the analyzer in the *worst case*. When programs use data structures that our grammar captures well, the analysis usually terminates without using this k limitation, and yields meaningful results.

 $\llbracket C \rrbracket : \widehat{D} \to \widehat{D}$ $\llbracket x := \texttt{new} \rrbracket \left(\mathcal{G}, R \right) = \left(\left\{ \left(s[x \to a], \, g[a \to \langle b, c \rangle \,, b \to \texttt{nil}, c \to \texttt{nil}] \right) \, | \, (s,g) \in \mathcal{G} \, \right\}, R \right) \, \texttt{new} \, a, b, c \to \texttt{nil}$ $\llbracket x := \operatorname{nil} \rrbracket (\mathcal{G}, R) = \left(\left\{ (s[x \to a], g[a \to \operatorname{nil}]) \mid (s, g) \in \mathcal{G} \right\}, R \right) \text{ new } a$ $\llbracket x := y \rrbracket (\mathcal{G}, R) = \text{when } y \in \text{dom}(s) \text{ for all } (s, g) \in \mathcal{G},$ $\left(\left\{\left(s[x \to s(y)], g\right) \mid (s, g) \in \mathcal{G}\right\}, R\right)$ $\llbracket x \text{->} 0 := y \rrbracket (\mathcal{G}, R) = \text{when unfold}(\mathcal{G}, R, x) = \mathcal{G}' \text{ and } \forall (s, g) \in \mathcal{G}'. y \in \text{dom}(s),$ $\left(\left\{\left(s,g[s(x)\rightarrow\langle s(y),c\rangle\right]\mid (s,g)\in\mathcal{G}',\ g(s(x))=\langle b,c\rangle\right\},R\right)$ $\llbracket x:=y\text{->}0 \rrbracket (\mathcal{G},R) = \text{when } \mathsf{unfold}(\mathcal{G},R,y) = \mathcal{G}',$ $\left(\left\{\left(s[x \rightarrow b], g\right) \mid (s, g) \in \mathcal{G}', \ g(s(y)) = \langle b, c \rangle\right\}, R\right)$ $\begin{bmatrix} C_1; C_2 \end{bmatrix} (\mathcal{G}, R) = \begin{bmatrix} C_2 \end{bmatrix} (\begin{bmatrix} C_1 \end{bmatrix} (\mathcal{G}, R))$ $\begin{bmatrix} \text{if } B \ C_1 \ C_2 \end{bmatrix} (\mathcal{G}, R) = \begin{bmatrix} C_1 \end{bmatrix} (\begin{bmatrix} B \end{bmatrix} (\mathcal{G}, R)) \sqcup \begin{bmatrix} C_2 \end{bmatrix} (\begin{bmatrix} ! \ B \end{bmatrix} (\mathcal{G}, R))$ $\llbracket \text{while } B \ C \rrbracket (\mathcal{G}, R) = \llbracket ! B \rrbracket \quad \text{lfix} \stackrel{\sqsubseteq}{=} \lambda A \colon \widehat{D}_k^{\nabla}. \text{ normalize}(A \sqcup (\mathcal{G}, R) \sqcup \llbracket C \rrbracket (\llbracket B \rrbracket A))$ $\llbracket C \rrbracket A = \top \text{ (other cases)}$ $[\![B]\!]:\widehat{D}\to\widehat{D}$ $\boxed{\llbracket x = y \rrbracket (\mathcal{G}, R)} = \text{when } \mathsf{split}(\mathsf{split}((\mathcal{G}, R), x), y) = (\mathcal{G}', R')$ $(\{(s,g) \in \mathcal{G}' \mid s(x) \equiv s(y) \lor g(s(x)) = g(s(y)) = \operatorname{nil})\}, R')$ $\llbracket ! x = y \rrbracket (\mathcal{G}, R) = \text{when split}(\text{split}((\mathcal{G}, R), x), y) = (\mathcal{G}', R')$ $(\{(s,g) \in \mathcal{G}' \mid s(x) \not\equiv s(y) \land (g(s(x)) \neq \operatorname{nil} \lor g(s(y)) \neq \operatorname{nil})\}, R')$ $\llbracket ! (!B) \rrbracket (\mathcal{G}, R) = \llbracket B \rrbracket (\mathcal{G}, R)$ $\llbracket B \rrbracket A = \top$ (other cases) Subroutine unfold unrolls the definition of a grammar: $\begin{aligned} \mathsf{unfold}((s,g),R,x) &= \begin{cases} \{(s,g[a \to \langle b,c \rangle, b \to v, c \to u]) | \langle v,u \rangle \in R(a)\} \text{ if } g(s(x)) = \alpha \land \operatorname{nil} \not\in R(\alpha) \\ \{(s,g)\} & \text{ if } g(s(x)) \text{ is a pair} \\ \top & \text{ otherwise} \end{cases} \\ \mathsf{unfold}(\mathcal{G},R,x) &= \begin{array}{c} \bigcup_{(s,g) \in \mathcal{G}} \mathsf{unfold}((s,g),R,x) & \text{ if } \forall (s,g) \in \mathcal{G}. \text{ unfold}((s,g),R,x) \neq \top \\ & \text{ otherwise} \end{cases} \end{aligned}$ Subroutine split((s, g), R, x) changes (s, g) to (s', g') s.t. s'(x) means nil iff g'(s'(x)) = nil. $\operatorname{split}((s,g), R, x) = \operatorname{if}(\exists \alpha. g(s(x)) = \alpha \land R(\alpha) \supseteq \{\operatorname{nil}\} \land R(\alpha) \neq \{\operatorname{nil}\})$ then $(\{(s, g[s(x) \rightarrow \text{nil}]), (s, g[s(x) \rightarrow \beta])\}, R[\beta \rightarrow R(\alpha) - \{\text{nil}\}])$ for fresh β else if $(\exists \alpha. g(s(x)) = \alpha \land R(\alpha) = \{\text{nil}\})$ then $(\{(s, g[s(x) \rightarrow \text{nil}])\}, R)$ else $(\{(s,g)\}, R)$ $\mathsf{split}(\mathcal{G}, R, x) = \begin{array}{c} \overset{()}{\vdash} {}_{(s,g)\in\mathcal{G}}\mathsf{split}((s,g), R, x) \text{ if } \forall (s,g)\in\mathcal{G}. x \in \mathrm{dom}(s) \\ \top & \mathrm{otherwise} \end{array}$ The algorithmic order \doteq defined in [8] satisfies that if $A \doteq B$, means(A) \Rightarrow means(B)

Fig. 3. Analysis.

7 Full Analysis

The basic version of our analysis, which we have presented so far, cannot deal with data structures with sharing, such as doubly linked lists and binomial heaps. If a program uses such data structures, the analysis gives up and returns \top .

The full analysis overcomes this shortcoming by using a more expressive language for a grammar, where a nonterminal is allowed to have parameters. The main feature of this new parameterized grammar is that an invariant for a data structure with sharing is expressible by a grammar, as long as the sharing is "cyclic." A parameter plays a role of "targets" of such cycles.

The overall structure of the full analysis is almost identical to the basic version in Figure 3. Only the subroutines, such as **normalize**, are modified. In this section, we will explain the full analysis by focusing on the new parameterized grammar, and the modified normalization function for this grammar. The full definition is in [8].

7.1 Abstract Domain

Let **self** and **arg** be two different symbolic locations. In the full analysis, the domains for shape graphs and grammars are modified as follows:

$$\begin{split} & NTermApp \stackrel{\Delta}{=} NonTerm \times (SymLoc + \bot) \quad NTermAppR \stackrel{\Delta}{=} NonTerm \times (\{\texttt{self}, \texttt{arg}\} + \bot) \\ & Graph \stackrel{\Delta}{=} (Vars \rightharpoonup_{fin} SymLoc) \times (SymLoc \rightharpoonup_{fin} \{\texttt{nil}\} + SymLoc^2 + NTermApp) \\ & Grammar \stackrel{\Delta}{=} NonTerm \rightharpoonup_{fin} \wp_{nf}(\{\texttt{nil}\} + (\{\texttt{nil}\} + \{\texttt{self}, \texttt{arg}\} + NTermAppR)^2) \end{split}$$

The main change in the new definitions is that all the nonterminals have parameters. All the uses of nonterminals in the old definitions are replaced by the applications of nonterminals, and the declarations of nonterminals in a grammar can use two symbolic locations self and arg, as opposed to none, which denote the implicit self parameter and the explicit parameter.⁸ For instance, a doubly-linked list is defined by dll ::= nil | $\langle \arg, dll(self) \rangle$. This grammar maintains the invariant that arg points to the previous cell. So, the first field of a node always points to the previous cell, and the second field the the next cell. Note that \perp can be applied to a nonterminal; this means that we consider subcases of the nonterminal where the arg parameter is not used. For instance, if a grammar R maps β to {nil, $\langle \arg, \arg \rangle$ }, then $\beta(\perp)$ excludes $\langle \arg, \arg \rangle$, and means the empty heap object.

As in the basic case, the precise meaning of a shape graph and a grammar is given by a translation into separation-logic assertions. We can define a translation means by modifying only means_v and means_g.

$$\begin{array}{ll} \operatorname{\mathsf{means}}_{\mathrm{v}}(a,\operatorname{nil}) \stackrel{\Delta}{=} a \doteq \operatorname{nil} & \operatorname{\mathsf{means}}_{\mathrm{v}}(a,\alpha(b)) \stackrel{\Delta}{=} \alpha(a,b) \\ \operatorname{\mathsf{means}}_{\mathrm{v}}(a,b) \stackrel{\Delta}{=} a \doteq b & \operatorname{\mathsf{means}}_{\mathrm{v}}(a,\alpha(\bot)) \stackrel{\Delta}{=} \forall b.\alpha(a,b) \\ \end{array}$$

In the last clause, b is a different variable from a. The meaning of a grammar is a context defining a set of recursive predicates.

$$\begin{array}{rcl} \operatorname{means}_{\mathrm{gc}}(R) & \stackrel{\Delta}{=} & \{\alpha(a,b) = \bigvee_{e \in R(\alpha)} \operatorname{means}_{\mathrm{gc}}(a,b,e)\}_{\alpha \in \operatorname{dom}(R)} \\ \operatorname{means}_{\mathrm{gc}}(a,b,\operatorname{nil}) & \stackrel{\Delta}{=} & \operatorname{means}_{\mathrm{v}}(a,\operatorname{nil}) \\ \operatorname{means}_{\mathrm{gc}}(a,b,\langle v_1,v_2\rangle) & \stackrel{\Delta}{=} & \exists a_1a_2.\,(a \mapsto a_1,a_2) \ast \operatorname{means}_{\mathrm{v}}(a_1,v_1\,\{a/\operatorname{self},b/\operatorname{arg}\}) \\ & \ast \operatorname{means}_{\mathrm{v}}(a_2,v_2\,\{a/\operatorname{self},b/\operatorname{arg}\}) \end{array}$$

In the second clause, a_1 and a_2 are variables that do not appear in v_1 , v_2 , a, b.

7.2 Normalization Function

To fully exploit the increased expressivity of the abstract domain, we change the normalization function in the full analysis. The most important change in the new normalization function is the addition of new rules (**cut**) and (**bfold**) into the **fold** procedure.

⁸ We allow only "one" explicit parameter. So, we can use pre-defined name arg.



Fig. 4. Examples of (cut) and (bfold).

The (\mathbf{cut}) rule enables the conversion of a cyclic structure to grammar definitions. Recall that the (\mathbf{fold}) rule can recognize a heap object only when the object does not have shared cells internally. The key idea is to "cut" a "non-critical" link to a shared cell, and represent the removed link by a parameter to a nonterminal. If enough such links are cut from an heap object, the object no longer has (explicitly) shared cells, so that the wrapping step of (**fold**) can be applied. The formal definition of the (**cut**) rule is:

(cut) $(\mathcal{G} \uplus \{(s, g \cdot [a \to \langle a_1, a_2 \rangle])\}, R) \rightsquigarrow \qquad \begin{array}{l} \mathcal{G} \cup \{(s, g \cdot [a \to \alpha(b)])\}, \\ R \cdot [\alpha \to \{\langle a_1, a_2 \rangle \{\mathsf{self}/a, \mathsf{arg}/b\}\}] \\ \text{where there are paths from variables to } a_1 \text{ and } a_2 \text{ in } g, \operatorname{free}(\langle v_1, v_2 \rangle) \subseteq \{a, b\}, \\ \text{and } \alpha \text{ is fresh. (If free}(\langle v_1, v_2 \rangle) \subseteq \{a\}, \text{ we use } \alpha(\bot) \text{ instead of } \alpha(b).) \end{array}$

Figure 4.(a) shows how a cyclic structure is converted to grammar definitions.⁹ In the first shape graph, "cell" a is shared because variable x points to a and "cell" c points to a, but the link from c to a is not critical because even without this link, a is still reachable from x. Thus, the (**cut**) rule cuts the link from c to a, introduces a nonterminal α_c with the definition $\{\langle \arg \rangle\}$, and annotates node c with $\alpha_c(a)$. Note that the resulting graph (the second shape graph in Figure 4.(a)) does not have explicit sharing. So, we can apply the (**fold**) rule to c, and then to b as shown in the last two shape graphs in Figure 4.(a).

The (**bfold**) rule wraps a cell "from the back." Recall that the (**fold**) rule puts a cell at the front of a heap object; it adds the cell as a root for a nonterminal. The (**bfold**) rule, on the other hand, puts a cell a at the exit of a heap object. When b is used as a parameter for a nonterminal α , the rule "combines" b and α . This rule can best be explained using a list-traversing algorithm. Consider a program that traverses a linked list, where variable r points to the head cell of the list, and variable c to the current cell of the list. The usual loop invariant of such a program is expressed by the first shape graph in Figure 4.(b). However, only with the (**fold**) rule, which adds a cell to the front, we cannot discover this invariant; one iteration of the program moves c to the next cell, and thus changes the shape graph into the second shape graph in Figure 4.(b), but this new graph is not similar to the initial one. The (**bfold**) rule changes the new shape graph back to the one for the invariant, by merging $\alpha(b)$ with cell b. The (**cut**) rule

⁹ To simplify the presentation, we assume that each cell in the figure has only a single field.

first cuts the link from b to c, extends a grammar with $[\gamma \rightarrow \{\langle \arg \rangle\}]$, and annotates the node b with $\gamma(c)$. Then, the (**bfold**) rule finds all the places where **arg** is used as itself in the definition of α , and replaces **arg** there by $\gamma(\arg)$. Finally, the rule changes the binding for a from $\alpha(b)$ to $\alpha(c)$, and eliminates cell b, thus resulting the last shape graph in Figure 4(b).¹⁰ The precise definition of (**bfold**) does what we call *linearity* check, in order to ensure the soundness of replacing **arg** by nonterminals:¹¹

 $\begin{array}{l} \textbf{(bfold)} \ (\mathcal{G} \cup \{(s,g \cdot [a \rightarrow \alpha(b),b \rightarrow \beta(w)])\}, R) \rightsquigarrow (\mathcal{G} \cup \{(s,g \cdot [a \rightarrow \alpha'(w)])\}, R \cdot [\alpha' \rightarrow E]) \\ \text{where } b \text{ does not appear in } g, \alpha \text{ is linear (that is, arg appears exactly once in each case of } R(\alpha)), \text{ and } E = \{\text{nil} \in R(\alpha)\} \cup \{\langle f(v_1), f(v_2) \rangle \mid \langle v_1, v_2 \rangle \in R(\alpha)\} \\ \text{where } f(v) = \text{if } (v \equiv \text{arg) then } \beta(\text{arg}) \text{ else } \text{ if } (v \equiv \alpha(\text{arg})) \text{ then } \alpha'(\text{arg}) \text{else } v \end{array}$

7.3 Correctness

The correctness of our analysis is expressed by the following theorem:

Theorem 1. For all programs C and abstract states (\mathcal{G}, R) , if $\llbracket C \rrbracket (\mathcal{G}, R)$ is a non- \top abstract state (\mathcal{G}', R') , then triple {means (\mathcal{G}, R) }C{means (\mathcal{G}', R') } holds in separation logic.

We proved this theorem in two steps. First, we showed a lemma that all subroutines, such as normalize and unfold, and the backward analysis are correct. Then, with this lemma, we applied the induction on the structure of C, and showed that $\{\text{means}(\mathcal{G}, R)\}C\{\text{means}(\mathcal{G}', R')\}$ is derivable in separation logic. The validity of the triple now follows, because separation-logic proof rules are sound. The details are in [8].

8 Experiments

We have tested our analysis with the six programs in Table 1. For each of the programs, we ran the analyzer, and obtained abstract states for a loop invariant and the result. In this section, we will explain the cases of binomial heap construction and the Schorr-Waite tree traversal. The others are explained at http://ropas.snu.ac.kr/grammar.

Binomial Heap Construction In this experiment, we took an implementation of binomial heap construction in [1], where each cell has three pointers: one to the left-most child, another to the next sibling, and the third to the parent. We ran the analyzer with this binomial heap construction program and the empty abstract state $(\{\}, [])$. Then, the analyzer inferred the following same abstract state (\mathcal{G}, R) for the result of the construction as well as for the loop invariant. Here we omit \perp from forest(\perp).

¹⁰ The grammar is slightly different from the one for the invariant. However, if we combine two abstract states and apply unify and simplify, then the grammar for the invariant is recovered.

¹¹ Here we present only for the case that the parameter of α is not passed to another different nonterminals. With such nonterminals, we need to do a more serious linearity check on those nonterminals, before modifying the grammar.

program	description	$\cos(\sec)$	analysis result
listrev.c	list construction followed by list	0.01	the result is a list
	reversal		
dbinary.c	construction of a tree with par-	0.01	the result is a tree with parent
	ent pointers		pointers
dll.c	doubly-linked list construction	0.01	the result is a doubly-linked list
bh.c	binomial heap construction	0.14	the result is a binomial heap
sw.c	Schorr-Waite tree traversal	0.05	the result is a tree
swfree.c	Schorr-Waite tree disposal	0.02	the tree is completely disposed

For all the examples, our analyzer proves the absence of null pointer dereference errors and memory leaks.

 Table 1. Experimental Results

$$\mathcal{G} = [x \rightarrow a], [a \rightarrow \mathsf{forest}] \qquad R = \begin{array}{c} \mathsf{forest} ::= \mathrm{nil} \mid \langle \mathsf{stree}(\mathsf{self}), \mathsf{forest}, \mathrm{nil} \rangle, \\ \mathsf{stree} ::= \mathrm{nil} \mid \langle \mathsf{stree}(\mathsf{self}), \mathsf{stree}(\mathsf{arg}), \mathsf{arg} \rangle \end{array}$$

The unique shape graph in \mathcal{G} means that the heap has only a single heap object whose root is stored in x, and the heap object is an instance of forest. Grammar Rdefines the structure of this heap object. It says that the heap object is a linked list of instances of stree, and that each instance of stree in the list is given the address of the containing list cell. These instances of stree are, indeed, precisely those trees with pointers to the left-most children and to the next sibling, and the parent pointer.

Schorr-Waite Tree Traversal We used the following (\mathcal{G}_0, R_0) as an initial abstract state:

 $\mathcal{G}_0 = \{([x \to a], [a \to \mathsf{tree}])\} \quad R_0 = [\mathsf{tree} ::= \mathrm{nil} \mid \langle \mathsf{I}, \mathsf{tree}, \mathsf{tree} \rangle]$

Here we omit \perp from tree(\perp). This abstract state means that the initial heap contains only a binary tree *a* whose cells are marked I.

Given the traversing algorithm and the abstract state (\mathcal{G}_0, R_0) , the analyzer produced (\mathcal{G}_1, R_1) for final states, and (\mathcal{G}_2, R_2) for a loop invariant:

 $\mathcal{G}_1 = [x \rightarrow a], [a \rightarrow \mathsf{treeR}] \qquad \qquad R_1 = [\mathsf{treeR} ::= \mathrm{nil} \mid \langle \mathbf{R}, \mathsf{treeR}, \mathsf{treeR} \rangle]$

 $\mathcal{G}_2 = [x \rightarrow a, y \rightarrow b], [a \rightarrow \mathsf{treeRI}, b \rightarrow \mathsf{rtree}]$

 $\begin{array}{ll} R_2 = & \operatorname{rtree}::=\operatorname{nil}|\left< \mathrm{R},\operatorname{treeR},\operatorname{rtree}\right>|\left< \mathrm{L},\operatorname{rtree},\operatorname{tree}\right>, & \operatorname{tree}::=\operatorname{nil}|\left< \mathrm{I},\operatorname{tree},\operatorname{tree}\right>, \\ & \operatorname{treeR}::=\operatorname{nil}|\left< \mathrm{R},\operatorname{treeR},\operatorname{treeR}\right>, & \operatorname{treeRI}::=\operatorname{nil}|\left< \mathrm{I},\operatorname{tree},\operatorname{tree}\right>|\left< \mathrm{R},\operatorname{treeR},\operatorname{treeR}\right> \end{array}$

The abstract state (\mathcal{G}_1, R_1) means that the heap contains only a single heap object x, and that this heap object is a binary tree containing only R-marked cells. Note that this abstract state implies the absence of memory leaks because the tree x is the only thing in the heap.

The loop invariant (\mathcal{G}_2, R_2) means that the heap contains two disjoint heap objects x and y. Since the heap object x is an instance of treeRI, the object x is an I-marked binary tree, or an R-marked binary tree. This first case indicates that x is first visited, and the second case that x has been visited before. The nonterminal rtree for the other heap object y implies that one of left or right field of cell y is reversed. The second case, $\langle R, \text{treeR}, \text{rtree} \rangle$, in the definition of rtree means that the current cell is marked R, its right field is reversed, and the left subtree is an R-marked binary tree. The third case, $\langle L, \text{rtree}, \text{tree} \rangle$, means that the current cell is marked L, the left field is reversed,

and the right subtree is an I-marked binary tree. Note that this invariant, indeed, holds because y points to the parent of x, so the left or right field of cell y must be reversed.

References

- Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. Introduction to Algorithms. MIT Press and McGraw-Hill Book Company, 2001.
- Patrick Cousot and Radhia Cousot. Abstract interpretation: A unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Proceedings of the ACM Symposium on Principles of Programming Languages*, pages 238–252, January 1977.
- Patrick Cousot and Radhia Cousot. Abstract interpretation frameworks. J. Logic and Comput., 2(4):511–547, 1992.
- 4. Patrick Cousot and Radhia Cousot. Formal language, grammar and set-constraintbased program analysis by abstract interpretation. In *Proceedings of the ACM Conference on Functional Programming Languages and Computer Architecture*, pages 170–181, La Jolla, California, June 1995. ACM Press, New York, NY.
- A. Deutsch. Interprocedural may-alias analysis for pointers: Beyond k-limiting. In Proceedings of the ACM Conference on Programming Language Design and Implementation, pages 230-241. ACM Press, 1994.
- Pascal Fradet and Daniel Le Métayer. Shape types. In Proceedings of the ACM Symposium on Principles of Programming Languages, pages 27–39. ACM Press, January 1997.
- 7. Nils Klarlund and Michael I. Schwartzbach. Graph types. In *Proceedings of the* ACM Symposium on Principles of Programming Languages, January 1993.
- Oukseh Lee, Hongseok Yang, and Kwangkeun Yi. Automatic verification of pointer programs using grammar-based shape analysis. Tech. Memo. ROPAS-2005-23, Programming Research Laboratory, School of Computer Science & Engineering, Seoul National University, March 2005.
- R. Manevich, M. Sagiv, G. Ramalingam, and J. Field. Partially disjunctive heap abstraction. In *Proceedings of the International Symposium on Static Analysis*, volume 3148 of *Lecture Notes in Computer Science*, pages 265–279, August 2004.
- A. Møller and M. I. Schwartzbach. The pointer assertion logic engine. In Proceedings of the ACM Conference on Programming Language Design and Implementation. ACM, June 2001.
- Peter W. O'Hearn, Hongseok Yang, and John C. Reynolds. Separation and information hiding. In Proceedings of the ACM Symposium on Principles of Programming Languages, pages 268–280. ACM Press, January 2004.
- John C. Reynolds. Separation logic: A logic for shared mutable data structures. In Proceedings of the 17th IEEE Symposium on Logic in Computer Science, pages 55–74. IEEE, July 2002.
- M. Sagiv, T. Reps, and R. Wilhelm. Solving shape-analysis problems in languages with destructive updating. ACM Trans. Program. Lang. Syst., 20(1):1–50, January 1998.
- Mooly Sagiv, Thomas Reps, and Reinhard Wilhelm. Parametric shape analysis via 3-valued logic. ACM Trans. Program. Lang. Syst., 24(3):217–298, 2002.
- Élodie-Jane Sims. Extending separation logic with fixpoints and postponed substitution. In Proceedings of the International Conference on Algebraic Methodology and Software Technology, volume 3116 of Lecture Notes in Computer Science, pages 475–490, 2004.