



**HAL**  
open science

## A Product of Shape and Sequence Abstractions

Josselin Giet, Félix Ridoux, Xavier Rival

► **To cite this version:**

Josselin Giet, Félix Ridoux, Xavier Rival. A Product of Shape and Sequence Abstractions. Static Analysis: 30th International Symposium, SAS 2023, Oct 2023, Cascais, Portugal. hal-04253341

**HAL Id: hal-04253341**

**<https://hal.science/hal-04253341>**

Submitted on 22 Oct 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A Product of Shape and Sequence Abstractions

Josselin Giet<sup>1</sup>, Félix Ridoux<sup>2,3</sup>, and Xavier Rival<sup>1</sup>

<sup>1</sup> INRIA Paris/CNRS/École Normale Supérieure/PSL Research University

<sup>2</sup> IMDEA Software Institute, Madrid, Spain

<sup>3</sup> Univ Rennes, F-35000 Rennes, France

`firstname.lastname@{ens.fr|ens-rennes.fr}`

**Abstract.** Traditional separation logic-based shape analyses utilize inductive summarizing predicates so as to capture general properties of the layout of data-structures, to verify accurate manipulations of, e.g., various forms of lists or trees. However, they also usually abstract away contents properties, so that they may only verify memory safety and invariance of data-structure shapes. In this paper, we introduce a novel abstract domain to describe sequences of values of unbounded size, and track constraints on their length and on extremal values contained in them. We define a reduced product of such a sequence abstraction together with an existing shape abstraction so as to infer both shape and contents properties of data-structures. We report on the implementation of the sequence domain, its integration into a static analyzer for C code, and we evaluate its ability to verify partial functional correctness properties for list and tree algorithms.



## 1 Introduction

Dynamically allocated data-structures based on lists, trees or graphs are common due to their flexibility as containers. However, programs using them are notoriously difficult to get right, especially in presence of destructive updates. Indeed, the correctness of such programs relies on a wide spectrum of properties that comprise memory safety (the absence of illegal pointer operations such as the dereference of a null pointer), the preservation of structural invariants like acyclicity, and subtle functional properties and relationships between the structure layout and its contents such as sortedness. For instance, let us consider a program that inserts an element in a binary search tree. First, it should not cause any runtime error or memory leak. Second, it should not create a cycle or break the tree structure. Third, it should preserve the binary search tree property and be functionally correct, namely ensure that the elements in the tree after insertion are the same as before plus the new element, inserted at the correct position, with respect to the order.

```

    struct tree { struct tree *l, *r; int d; };

tree( $\alpha$ ) :=
  | emp  $\wedge \alpha = \mathbf{0x0}$ 
  |  $\exists \alpha_l, \alpha_r, \delta, \alpha. \mathbf{l} \mapsto \alpha_l * \alpha. \mathbf{r} \mapsto \alpha_r * \alpha. \mathbf{d} \mapsto \delta * \mathbf{tree}(\alpha_l) * \mathbf{tree}(\alpha_r) \wedge \alpha \neq \mathbf{0x0}$ 
trees( $\alpha, S$ ) :=
  | emp  $\wedge \alpha = \mathbf{0x0}$ 
  |  $\exists \alpha_l, \alpha_r, \delta, S_l, S_r, \alpha. \mathbf{l} \mapsto \alpha_l * \alpha. \mathbf{r} \mapsto \alpha_r * \alpha. \mathbf{d} \mapsto \delta * \mathbf{tree}_s(\alpha_l, S_l) * \mathbf{tree}_s(\alpha_r, S_r)$ 
     $\wedge \alpha \neq \mathbf{0x0} \wedge S = S_l. [\delta]. S_r$ 

```

**Fig. 1.** A C tree data-type and associated inductive summarizing predicates

Abstract interpretation [18] provides a general framework to build a sound static analysis from a basic semantics and an abstraction relation, and to verify semantic properties. Notably, it has been applied to verify numerical properties [21,36], the absence of runtime errors [7], string properties [26,4], array properties [29,30,31,20], liveness properties [57], and security properties [5,27,25]. Several families of shape analyses have also been designed to infer properties of programs manipulating dynamic data-structures, including TVLA [53] and shape analyses based on separation logic [52]. They can reason over structures like lists [13,14] or more general families of structures with an inductive layout [15,32] such as binary trees.

However, few shape analyses reason not only about the layout of data-structures but also about their contents, so as to verify, e.g., that a container consists of the expected collection of elements with the expected multiplicity. While [41,22] handle set predicates, they do not track properties related to element multiplicities or order. Similarly, [43] handles sorting properties of specific families of composite structures in arrays but does not consider general lists or trees. The analyses presented in [10,9,12] precisely abstract singly-linked lists storing numerical data. They compute numerical properties over these data such as "variable  $x$  is the sum of all elements in list  $l$ ", or relation between element values and indexes to express sorting. However, it does not handle trees or doubly linked lists. Therefore, in this paper, we seek for an abstraction of data-structure contents that can verify complex invariants (e.g., involving elements orders or multiplicities) as well as some functional properties (like sorting). To illustrate our approach, we consider the classical tree type definition shown in Figure 1 and assume that we only consider acyclic instances. The inductive predicate **tree** summarizes valid memory regions storing exactly a complete and acyclic tree. More precisely, the predicate **tree**( $\alpha$ ) either describes an empty tree (then,  $\alpha$  is the null pointer), or a memory region where  $\alpha$  points to a valid **tree** block, the **l** and **r** fields of which point to the roots of disjoint (possibly empty) subtrees, as expressed by predicates **tree**( $\alpha_l$ ) and **tree**( $\alpha_r$ ). Note that separating conjunction  $*$  [52] combines disjoint memory regions. A basic region is either an atomic cell described by a points-to predicate such as  $\alpha. \mathbf{d} \mapsto \delta$  or an instance of some inductive predicate. While predicate **tree** describes the layout of memory cells and pointers, it does not convey any information about their contents. By contrast, **tree**<sub>s</sub> extends **tree** with an additional symbolic parameter  $S$  to expose the sequence of values stored in the tree, read from left to

right. When the tree is empty, so is its sequence of elements. The sequence stored in a non-empty tree is obtained by first considering the left subtree, then the contents of the root node and finally the right tree. If we additionally require the elements of  $S$  be sorted, then  $\mathbf{tree}_s(\alpha, S)$  describes binary search trees with root  $\alpha$ .

An advantage of this approach is that it allows to split the abstraction into two rather independent components, namely a separation logic based abstraction of the data-structures and another abstraction for properties of sequences of values stored in them. While [41] extends inductive predicates in a similar manner, it only supports set constraints. Therefore, we introduce a new abstract domain devoted to the representation of constraints over sequences. Existing sequence abstractions typically rely on regular expressions or finite automata [47,3,49]. More recently, [4] extends such an abstraction with sub-string, length, and element position constraints. However, these abstractions lack predicates such as constraints over extremal elements or sortedness. Our sequence abstract domain expresses not only relational constraints (it can express that a symbolic sequence is a fragment of another) but also constraints over length, extremal values, and specific predicates like sortedness. Although we use this sequence abstract domain for shape analysis, it could be used independently for other kinds of analyses.

To take advantage of this abstraction in shape analysis, we define a reduced product with a separation logic-based shape abstract domain. This product ties symbolic parameters of inductive predicates in separation logic together with sequence constraints. Sequence constraints that are inferred during the analysis (for instance when unfolding inductive predicates) are passed to the sequence domain. The reduced product also ensures communication between both domains for the computation of abstract operators such as union.

To summarize, we make the following contributions:

- After we overview our analysis in Section 2, we introduce a relational abstract domain dedicated to reasoning over sequences in Section 3;
- We define a reduced product between the new sequence domain and a separation logic-based abstract domain so as to extend a shape analysis with sequence reasoning capability. We first introduce the basic elements of the reduced product in Section 4 in the context of singly-linked lists. We discuss issues related to general inductive predicates in Section 5.
- We report on the implementation of our analysis in the MemCAD static analyzer [40] and on its evaluation in Section 6. We show that it can cope with the verification of sorting programs and operations over binary search trees.

## 2 Overview

In this section, we give an overview of the main principles of our static analysis by demonstrating it on the insertion program shown in Figure 2. When applied to a binary search tree, this function inserts an element at the expected position to preserve sortedness. We study the verification of functional correctness expressed as partial correctness with respect to a pre-condition and a post-condition (Figure 2). To formalize these, we let  $\mathbf{sort}$  be a symbolic function over sequences of values that

```

1  // assume tree_s(t, S) ∧ S = sort(S)
2  if(t == null){
3  // ...
4  }else{
5  struct tree* c = t;
6  while(c->d <= i && c->l != null ||
7  c->d > i && c->r != null)
8  c = (c->d <= i) ? c->l : c->r;
9  // ...
10 }
11 } // assert tree_s(t, sort(S.[i]))

```

Fig. 2. Function for insertion in a binary search tree

$$(\&t \mapsto \alpha_0 * \&c \mapsto \alpha_0 * \mathbf{tree}_s(\alpha_0, S)) \quad \wedge \quad (S = \mathbf{sort}(S) \wedge \alpha_0 \neq \mathbf{0x0})$$

(a) Abstract state at the end of line 6

$$\left( \begin{array}{l} \&t \mapsto \alpha_0 * \&c \mapsto \alpha_1 \\ * \alpha_0.l \mapsto \alpha_1 * \mathbf{tree}_s(\alpha_1, S_l) \\ * \alpha_0.d \mapsto \delta \\ * \alpha_0.r \mapsto \alpha_2 * \mathbf{tree}_s(\alpha_2, S_r) \end{array} \right) \quad \wedge \quad \left( \begin{array}{l} S = S_l.[\delta].S_r \wedge S = \mathbf{sort}(S) \\ \wedge S_l = \mathbf{sort}(S_l) \wedge S_r = \mathbf{sort}(S_r) \\ \wedge \max_{S_l} \leq \delta \leq \max_{S_r} \\ \wedge \delta \leq i \wedge \alpha_0, \alpha_1 \neq \mathbf{0x0} \end{array} \right)$$

(b) Abstract state at the end of line 9, first case of the condition

$$\left( \begin{array}{l} \&t \mapsto \alpha_0 * \&c \mapsto \alpha' \\ * \mathbf{treeseq}_s(\alpha_0, \alpha', S_1 \sqcup S_2) \\ * \mathbf{tree}_s(\alpha', S_0) \end{array} \right) \quad \wedge \quad \left( \begin{array}{l} S = S_1.S_0.S_2 \wedge S = \mathbf{sort}(S) \\ \wedge S_i = \mathbf{sort}(S_i) \quad i \in \{0, 1, 2\} \\ \wedge \max_{S_1} \leq i \leq \max_{S_2} \wedge \max_{S_1} \leq \max_{S_0} \\ \wedge \min_{S_0} \leq \max_{S_2} \wedge \alpha_0, \alpha' \neq \mathbf{0x0} \end{array} \right)$$

(c) Abstract state after the first widening

Fig. 3. Selected abstract states

maps any sequence to its sorted permutation. Then, the pre-condition assumption assumes that  $\mathbf{t}$  is a well-formed tree described with predicate  $\mathbf{tree}_s(\mathbf{t}, S)$  and such that  $S = \mathbf{sort}(S)$  (i.e., such that the elements  $S$  in  $\mathbf{t}$  are sorted). Likewise, the post-condition asserts that  $\mathbf{t}$  is still a well-formed tree, the contents of which is sorted and comprises exactly the elements in  $S$  plus the added value  $i$ .

We now discuss the abstraction used by our static analysis. We combine an existing memory abstraction, inspired by separation logic-based shape analyses such as [13,15], a relational numerical abstraction such as convex polyhedra [21], and a novel abstract domain for sequences. Intuitively, the latter describes conjunctions of constraints over both symbolic sequences of values (such as  $S$ ) and values manipulated by the program. These constraints consist of equalities of pairs of symbolic sequence expressions such as  $S' = \mathbf{sort}(S.[i])$ . Moreover, the inductive predicates used in the memory abstraction are instances of the  $\mathbf{tree}_s$  predicate of Figure 1. For instance, the abstract pre-condition simply consists of the memory predicate  $\mathbf{tree}_s(\mathbf{t}, S)$  and the sequence predicate  $S = \mathbf{sort}(S)$ , for some existentially quantified symbolic sequence variable  $S$ .

The analysis proceeds by forward abstract interpretation [18]: it computes over-approximate abstract post-conditions for basic statements, and uses widening to enforce the convergence of abstract iterations for loops. Since the analysis uses a reduced product [19], an abstract state consists of a pair of components, namely the shape abstraction that describes the layout of data-structures and the contents’ abstraction made of constraints over values and sequences of values. For each analysis step, information stored in either component may be used in order to refine the other, which we discuss next.

We focus on the analysis of the loop that searches for the insertion point in the **else** branch. First, the analysis of the condition test enriches the pre-condition with the constraint  $\mathfrak{t} \neq \mathbf{null}$  as shown in the abstract state in Figure 3(a). Then, the analysis continues with the loop. The condition is a disjunction thus the analysis considers each case separately. For the first case, it refines the abstract state to reflect that the condition  $\mathfrak{c} \rightarrow \mathfrak{d} \leq \mathfrak{i} \ \&\& \ \mathfrak{c} \rightarrow \mathfrak{l} \neq \mathbf{null}$  evaluates to **true**. Since both memory cells  $\mathfrak{c} \rightarrow \mathfrak{d}$  and  $\mathfrak{c} \rightarrow \mathfrak{l}$  are abstracted by the predicate  $\mathbf{tree}_s(\alpha_0, S)$ , this predicate needs to be unfolded to enable the analysis of the condition. The first disjunct of inductive predicate  $\mathbf{tree}_s$  (Figure 1), which corresponds to the null pointer, is ruled out by constraint  $\mathfrak{t} \neq \mathbf{null}$ . Therefore, only the second disjunct (non-empty tree) needs be considered. This shows how one component of the abstract state can refine the other. Thus, the analysis generates a new abstract state that exposes the root of the tree and lets  $\alpha_1$ ,  $\alpha_2$ , and  $\delta$  denote the contents of its  $\mathfrak{l}$ ,  $\mathfrak{r}$ , and  $\mathfrak{d}$  fields. We remark that the inductive predicate unfolding also splits the symbolic sequence into  $S = S_l. [\delta]. S_r$ . Then the sequence domain derives that  $S_l$  and  $S_r$  are sorted since they are subsequences of a sorted sequence. It also infers that all values in  $S_l$  are less than  $\delta$  that is itself less than all values in  $S_r$  by definition of **sort**, which writes down  $\max_{S_l} \leq \delta \leq \min_{S_r}$ . Last, it also retains the numerical constraint  $\mathfrak{i} \leq \delta$ . Figure 3(b) shows the resulting abstract state after the assignment line 9. In the case of the other disjunct, the tree is also unfolded but  $\mathfrak{c}$  points to  $\alpha_2$  instead of  $\alpha_1$  and the constraint over  $\mathfrak{i}$  and  $\delta$  is  $\delta < \mathfrak{i}$ .

The widening operator over-approximates abstract union of successive abstract iterates at loop head. In this case, it generalizes abstract states such as the ones shown in Figure 3(a) and Figure 3(b) by weakening them locally. Indeed, in all three states,  $\mathfrak{c}$  points to a well-formed tree containing a sequence  $S_0$ . Moreover, the remaining of the memory region corresponds to a (possibly empty) partial tree: if it was completed by a tree with root pointed by  $\mathfrak{c}$ , the whole region would form a complete tree with root pointed by  $\mathfrak{t}$ . We call such a partial tree a *tree segment predicate* (the name *segments* comes from the analogy with list segments) and observe that it can be automatically derived from  $\mathbf{tree}_s$  and defined by induction in Figure 4. When widening synthesizes an instance of  $\mathbf{treeseg}_s$  in Figure 3(c), it needs to infer its sequence argument. The sequence  $S$  of elements stored into the whole tree can be split into three parts,  $S_0$ ,  $S_1$ , and  $S_2$  where  $S_0$  is the sequence of elements stored in the subtree pointed to by  $\mathfrak{c}$  and  $S_1$  (resp.,  $S_2$ ) denote the sequence of elements stored in the “left” (resp., “right”) part of the tree segment. This implies that the sequence argument of  $\mathbf{treeseg}_s$  is not a contiguous sequence. Therefore, it is represented as  $S_1 \sqcap S_2$  in the loop invariant Figure 3(c) where the

$$\begin{aligned}
\mathbf{treeseg}_s(\alpha, \alpha', S \sqsupset S') := & \\
| \mathbf{emp} \wedge \alpha = \alpha' \wedge S = S' = \square & \\
| \exists \alpha_l, \alpha_r, \mathbf{v}, S_l, S'_l, S_r, \alpha.l \mapsto \alpha_l * \alpha.r \mapsto \alpha_r * \alpha.d \mapsto \mathbf{v} * \mathbf{treeseg}_s(\alpha_l, \alpha', S_l \sqsupset S'_l) & \\
| * \mathbf{tree}_s(\alpha_r, S_r) \wedge \alpha \neq \mathbf{0x0} \wedge S = S_l \wedge S' = S'_l.[\mathbf{v}].S_r & \\
| \exists \alpha_l, \alpha_r, \mathbf{v}, S_l, S_r, S'_r, \alpha.l \mapsto \alpha_l * \alpha.r \mapsto \alpha_r * \alpha.d \mapsto \mathbf{v} * \mathbf{tree}_s(\alpha_l, S_l) & \\
| * \mathbf{treeseg}_s(\alpha_r, \alpha', S_r \sqsupset S'_r) \wedge \alpha \neq \mathbf{0x0} \wedge S = S_l.[\mathbf{v}].S_r \wedge S = S'_r &
\end{aligned}$$

**Fig. 4.** Inductive summarizing predicate describing tree segment

*placeholder* notation  $\sqsupset$  stands for the sequence of elements in the “missing subtree” of the segment. When composing  $\mathbf{treeseg}_s$  and  $\mathbf{tree}_s$  the analysis operations resolve sequences using such  $\sqsupset$  symbol. Based on this loop invariant, the analysis of the final few assignments of the insertion function produces an abstract state that implies the desired post-condition.

### 3 Abstract Domain for Sequences

In this section, we define the sequence abstract domain, including its elements and the constraints they denote, its concretization, and its main abstract operators.

#### 3.1 Sequences abstraction

An element of the abstract domain of sequences is a conjunction of constraints over a finite set of symbolic variables that stand either for sequences of base values, for base values, or for sets of values. Beside sequence equalities and predicates like sortedness, we also consider numerical upper/lower bounds over the values in sequences and multi-set constraints over the collections of values in sequences.

*Concrete states.* Let  $\mathbb{V}$  denote a set of values. Although  $\mathbb{V}$  usually denotes a set of scalar values (including addresses), our only assumptions on  $\mathbb{V}$  is that it has a total ordering  $\preceq$  with extremal values  $+\infty, -\infty$ . Since our domain constrains both variables that range over  $\mathbb{V}$  and variables that range over sequences of values in  $\mathbb{V}$ , we need several kinds of *symbolic variables*. In the following, we let symbols  $\alpha, \alpha_0, \alpha'_0, \beta, \dots \in \mathbb{X}_n$  denote *value symbolic variables*, namely, variables that stand for a value in  $\mathbb{V}$ . To express constraints on the set  $\mathbb{V}^*$  of all the finite words on alphabet  $\mathbb{V}$ , we let a separate set  $\mathbb{X}_s$ , represent *sequence symbolic variables*. We note  $S, S_1, S', P, \dots \in \mathbb{X}_s$  such sequence variables. Finally, we write  $\mathcal{M}(\mathbb{V})$  for the set of multisets of values in  $\mathbb{V}$  and let  $\mathbb{X}_m$  be the set of *multi-set valued symbolic variables*. Moreover, if  $S \in \mathbb{X}_s$  is a sequence variable, we attach to it three numerical variables  $\mathbf{len}_S, \mathbf{min}_S, \mathbf{max}_S$  in  $\mathbb{X}_n$  that respectively denote the length, minimum and maximum value of  $S$ , and that there exists a multi-set variable  $\mathbf{multi}_S \in \mathbb{X}_m$  that denotes the multi-set of its elements.

A *concrete state* comprises three functions that map each kind of symbolic variables to elements of the corresponding type. Due to the relationship between a

$$\begin{aligned}
 E & ::= \square \mid [\alpha] \mid S \mid E.E \mid \mathbf{sort}(E) & C \in \mathbb{C} & ::= S = E \\
 & \text{(a) Syntax of expressions } (E) \text{ and constraints } (C) \\
 \llbracket \square \rrbracket_s(\sigma) & = \varepsilon & \llbracket [\alpha] \rrbracket_s(\sigma) & = \sigma_n(\alpha) \\
 \llbracket S \rrbracket_s(\sigma) & = \sigma_s(S) & \llbracket E_1.E_2 \rrbracket_s(\sigma) & = \llbracket E_1 \rrbracket_s(\sigma). \llbracket E_2 \rrbracket_s(\sigma) \\
 \llbracket \mathbf{sort}(E) \rrbracket_s(\sigma) & = a_{\pi(1)} \dots a_{\pi(n)} \text{ where } & \begin{cases} \llbracket E \rrbracket(\sigma) = a_1 \dots a_n \\ \forall i \in [1, n-1], a_{\pi(i)} \preceq a_{\pi(i+1)} \\ \pi \text{ is a permutation of } [1, n] \end{cases} \\
 (\sigma_n, \sigma_s) \models_s S = E & \text{ iff } \sigma_s(S) = \llbracket E \rrbracket_s(\sigma_n, \sigma_s) \\
 & \text{(b) Semantics}
 \end{aligned}$$

**Fig. 5.** Sequence expressions and constraints: syntax and semantics

sequence symbolic variable  $S$  and  $\mathbf{len}_S$ ,  $\mathbf{min}_S$ ,  $\mathbf{max}_S$ , and  $\mathbf{multi}_S$ , a concrete state is valid if and only if it maps these five variables into consistent objects. Formally:

**Definition 1.** A concrete state is a tuple  $\sigma = (\sigma_n, \sigma_m, \sigma_s)$  where the functions  $\sigma_n : \mathbb{X}_n \rightarrow \mathbb{V}$ ,  $\sigma_m : \mathbb{X}_m \rightarrow \mathcal{M}(\mathbb{V})$ ,  $\sigma_s : \mathbb{X}_s \rightarrow \mathbb{V}^*$  are such that, for all  $S$  in  $\mathbb{X}_s$ ,

$$\begin{aligned}
 \sigma_s(S) = a_1 \dots a_n & \Rightarrow \begin{cases} \sigma_n(\mathbf{min}_S) = \min_i a_i \wedge \sigma_n(\mathbf{max}_S) = \max_i a_i \\ \wedge \sigma_n(\mathbf{len}_S) = n \wedge \sigma_m(\mathbf{multi}_S) = \{a_1, \dots, a_n\} \end{cases} \\
 \sigma_s(S) = \varepsilon & \Rightarrow \begin{cases} \sigma_n(\mathbf{min}_S) = +\infty \wedge \sigma_n(\mathbf{max}_S) = -\infty \\ \wedge \sigma_n(\mathbf{len}_S) = 0 \wedge \sigma_m(\mathbf{multi}_S) = \emptyset \end{cases}
 \end{aligned}$$

For short, given a state  $\sigma$ , we note its components  $\sigma_n$ ,  $\sigma_m$ , and  $\sigma_s$ . We write  $\Sigma$  for the set of all such concrete states.

*Abstract sequence constraints.* The sequence abstract domain relies on expressions and constraints over symbolic variables. Their syntax is shown in Figure 5(a). An expression is either the empty sequence, or a sequence of length one that consists of a value symbolic variable, or a sequence symbolic variable, or a concatenation of expressions, or the sorting of a sequence expression returned by the function  $\mathbf{sort} : E \rightarrow E$ , (introduced in Section 2). Given a state  $\sigma$ , a sequence expression  $E$  evaluates into a sequence of values  $\llbracket E \rrbracket_s(\sigma)$ , as shown in Figure 5(b). Sequence constraints are *definition constraints* of the form  $S = E$ , as shown in Figure 5(a). Allowing only symbolic sequence variables in the left-hand side of equalities somewhat limits expressiveness but simplifies the machine representation of abstract elements. The semantics of constraints is defined based on a satisfaction relation  $\models_s$  that is spelled out in Figure 5(b): we write  $(\sigma_n, \sigma_s) \models_s C$  when constraint  $C$  holds in concrete state  $(\sigma_n, \sigma_s)$ . We note  $\mathbb{C}$  for the set of sequence constraints.

*Parameter abstract domains.* In the following, we assume two abstract domains are fixed, taken as parameters by the sequence abstraction. First,  $\mathbb{D}_n^\sharp$  represents numerical constraints and provides a concretization function  $\gamma_n : \mathbb{D}_n^\sharp \rightarrow \mathcal{P}(\mathbb{X}_n \rightarrow \mathbb{V})$ . Possible choices for  $\mathbb{D}_n^\sharp$  include intervals [18], octagons [48], or convex polyhedra [21] abstract domains. Second,  $\mathbb{D}_m^\sharp$  represents multi-set constraints and provides a concretization function  $\gamma_m : \mathbb{D}_m^\sharp \rightarrow \mathcal{P}(\mathbb{X}_m \rightarrow \mathcal{M}(\mathbb{V}))$ . Our implementation uses a variation of the set domain of [41] that describes multi-set constraints.



*Sequence abstraction.* An abstract state consists either of a special element  $\perp$  that denotes the empty set of concrete states or of a finite conjunction of sequence constraints together with numerical and multi-set constraints:

**Definition 2 (Sequence abstraction).** *The abstract sequence domain  $\Sigma^\sharp$  is defined as  $\{\perp\} \uplus \{(\sigma_n^\sharp, \sigma_m^\sharp, C_0 \wedge \dots \wedge C_n) \mid \sigma_n^\sharp \in \mathbb{D}_n^\sharp, \sigma_m^\sharp \in \mathbb{D}_m^\sharp, C_0, \dots, C_n \in \mathbb{C}\}$ . Furthermore, its concretization  $\gamma_\Sigma : \Sigma^\sharp \rightarrow \mathcal{P}(\Sigma)$  is defined by  $\gamma_\Sigma(\perp) = \emptyset$  and:*

$$\gamma_\Sigma(\sigma_n^\sharp, \sigma_m^\sharp, C_0 \wedge \dots \wedge C_n) = \{(\sigma_n, \sigma_m, \sigma_s) \mid \sigma_n \in \gamma_n(\sigma_n^\sharp) \wedge \sigma_m \in \gamma_m(\sigma_m^\sharp) \wedge \forall i, \sigma_n, \sigma_s \models_s C_i\}$$

For consistency, we use  $\sigma_s^\sharp$  as a generic notation for a finite conjunction of constraints  $C_0 \wedge \dots \wedge C_n$  and  $\sigma^\sharp$  for a generic triple  $(\sigma_n^\sharp, \sigma_m^\sharp, \sigma_s^\sharp)$ . We remark that the empty conjunction of constraints concretizes into  $\Sigma$  thus we note it  $\top$ .

*Machine representation.* For the sake of algorithmic efficiency, we rely on an optimized machine representation for sequence constraints in non-bottom abstract states. First, we let equality constraints between variables be described by union-find data-structures, which enables the incremental computation of equality classes representatives. Emptiness constraints ( $S = []$ ) and sortedness constraints ( $S = \text{sort}(S)$ ) are marked by tags over sequence variables. Finally, other equality constraints are represented with a map data type, the keys of which are the left hand side variables. For instance,  $S = [\alpha]$  boils down to a map entry  $S \mapsto [\alpha]$ .

### 3.2 Abstract operations

We now discuss abstract operations on sequence abstract states. In this subsection, we discuss two operations: **guard** $_\Sigma$  refines an abstract sequence element into its conjunction with an additional constraint and **verify** $_\Sigma$  attempts to discharge a sequence constraint (so as to, e.g., verify an assertion). We assume that the underlying domains also implement similar operators. For instance, we require the numerical domain to provide an operator **guard** $_n$  that inputs a numerical constraint and a  $\sigma_n^\sharp \in \mathbb{D}_n^\sharp$  and refines the latter with that constraint.

*Abstract sequence condition.* First, we consider the *abstract sequence condition* operator **guard** $_\Sigma : \mathbb{C} \times \Sigma^\sharp \rightarrow \Sigma^\sharp$  which refines an abstract state with an additional sequence constraint. While a naive implementation of **guard** $_\Sigma(C, \sigma^\sharp)$  would simply add the constraint  $C$  to the conjunction  $\sigma_s^\sharp$  component, this would be imprecise in general. Indeed, the conjunction  $C \wedge \sigma_s^\sharp$  may be equivalent to  $\perp$ . Moreover,  $C \wedge \sigma_s^\sharp$  may entail constraints that are strictly more precise than those in  $\sigma_s^\sharp$ .

At a high level, **guard** $_\Sigma$  performs three kinds of operations:

1. *Compaction* simplifies constraints by rewriting the right hand side of definition constraints into the left hand side, wherever possible. For example,  $S = S'. [\alpha]. S'' \wedge S_1 = S'. [\alpha]$  simplifies into  $S = S_1. S'' \wedge S_1 = S'. [\alpha]$ .
2. *Saturation* synthesizes additional numerical and multi-set constraints that can be derived from a newly added constraint. For instance,  $S = S'. [\alpha]$  entails

that  $\mathbf{len}_S = 1 + \mathbf{len}_{S'}$ . Likewise, some constraints may entail that a sequence is empty. Another special kind of saturation occurs when the whole state can be reduced to  $\perp$  as incompatible constraints are detected. As saturation is the most complex part of  $\mathbf{guard}_{\Sigma}$ , we detail it below.

3. *Detection of cyclic constraints* prevents compaction and saturation from adding too many, redundant constraints, and it ensures the termination of algorithms iterating on definitions. We discuss this in Example 2.

We now discuss constraint saturation more in detail:

- The *length constraints saturation* derives numerical constraints from the equality of the length of both sides of a new definition constraint  $S = E$ . Indeed, such a constraint implies  $\mathbf{len}_S = \tau_{\mathbf{len}}(E)$ , which can be added to the  $\sigma_n^\sharp$  component using  $\mathbf{guard}_n$ , where  $\tau_{\mathbf{len}}$  is defined by:

$$\begin{array}{l} \tau_{\mathbf{len}}(\square) = 0 \quad \tau_{\mathbf{len}}(E.E') = \tau_{\mathbf{len}}(E) + \tau_{\mathbf{len}}(E') \quad \tau_{\mathbf{len}}(S) = \mathbf{len}_S \\ \tau_{\mathbf{len}}([\alpha]) = 1 \quad \tau_{\mathbf{len}}(\mathbf{sort}(E)) = \tau_{\mathbf{len}}(E) \end{array}$$

- The *multi-set contents constraints saturation* operates similarly, and derives multi-set equalities from definition constraints. Surely,  $S = E$  entails  $\mathbf{multi}_S = \tau_{\mathbf{mul}}(E)$ , which can refine the  $\sigma_m^\sharp$  part using  $\mathbf{guard}_m$  where  $\tau_{\mathbf{mul}}$  is defined by:

$$\begin{array}{l} \tau_{\mathbf{mul}}(\square) = \emptyset \quad \tau_{\mathbf{mul}}(E.E') = \tau_{\mathbf{mul}}(E) \uplus \tau_{\mathbf{mul}}(E') \quad \tau_{\mathbf{mul}}(S) = \mathbf{multi}_S \\ \tau_{\mathbf{mul}}([\alpha]) = \{\alpha\} \quad \tau_{\mathbf{mul}}(\mathbf{sort}(E)) = \tau_{\mathbf{mul}}(E) \end{array}$$

- The *detection of empty sequence variables* derives new definition constraints of the form  $S = \square$  when either sequence constraints or numerical constraints entail the emptiness of  $S$ . For instance:
  - when  $\sigma_s^\sharp$  contains constraints  $S' = \square$  and  $S'' = \square$ , the constraint  $S = S'.S''$  simplifies into  $S = \square$ ;
  - when  $\sigma_n^\sharp$  contains the constraint  $\mathbf{len}_S = 0$ , then it follows that  $S = \square$ .
- The *detection of sorted sequence variables* do the same for constraints of the form  $S = \mathbf{sort}(S)$  thanks to definitions of  $S$  and to numerical inequalities:

$$\frac{S = S_1 \dots S_n \quad \forall i, S_i = \mathbf{sort}(S_i) \quad \forall i < j, \mathbf{max}_{S_i} \leq \mathbf{max}_{S_j}}{S = \mathbf{sort}(S)}$$

Such rule is very costly as it checks a quadratic amount of numerical inequalities. Nevertheless, relaxing the rule by only considering the case  $j = i + 1$  is not sound, since sequence variables may be empty. Therefore, two consecutive elements in  $\sigma_s(S)$  can come from non-consecutive sequence variables.

- The *extremal values inequalities saturation* derives numerical inequalities from a definition constraint  $S = E$  by case analysis over the right hand side  $E$ , and can be summarized by a set of derivation rules. The rules below describe such reasoning steps:

$$\begin{array}{l} \frac{S = E \quad \alpha \in \mathbf{fv}(E)}{\mathbf{min}_S \leq \alpha \leq \mathbf{max}_S} \\ \frac{S = \square \quad \alpha \in \mathbb{X}_n}{\mathbf{max}_S < \alpha < \mathbf{min}_S} \end{array} \quad \frac{\frac{S = E \quad S' \in \mathbf{fv}(E)}{\mathbf{min}_S \leq \mathbf{min}_{S'} \quad \mathbf{max}_{S'} \leq \mathbf{max}_S}}{S' = \mathbf{sort}(S') \quad S' = \dots[\alpha].S \dots} \quad \frac{}{\alpha \leq \mathbf{min}_S}$$

- As an example, the first rule states that numerical constraints can be derived from the knowledge that  $S$  is a concatenation of several components including a numerical variable  $\alpha$ ; in this case novel numeric constraints expressing that  $\alpha$  is bounded by the extremal values of  $S$  can be added to  $\sigma_n^\#$  using operator  $\mathbf{guar}\mathbf{d}_n$ . Similarly, the second rule states that the extremal values of a sequence are bounded by the extremal values of any sequence containing it. The third rule states that an empty sequence supports arbitrary bounds. Finally, the fourth rule allows to reason over bounds when a sequence is known to be sorted.
- The *decomposition of equality constraints* synthesizes additional equality constraints that can be derived when two definition constraints  $S = E_0$  and  $S = E_1$  over the same name can be found in  $\sigma_s^\#$ . Indeed, when both  $E_0$  and  $E_1$  can be decomposed simultaneously, new equalities can be immediately derived:

$$\frac{[\alpha_0].E_0 = [\alpha_1].E_1}{\alpha_0 = \alpha_1 \quad E_0 = E_1} \quad \frac{S.E_0 = E_1 \quad S = []}{E_0 = E_1}$$

In less obvious decomposition cases, further constraints can still be derived with the help of the numerical constraints. Indeed:

$$\frac{S_0.E_0 = S_1.E_1 \quad \mathbf{len}_{S_0} = \mathbf{len}_{S_1}}{S_0 = S_1 \quad E_0 = E_1}$$

Obviously, this inference may take place only when  $\mathbf{len}_{S_0} = \mathbf{len}_{S_1}$  can be proved in the numerical domain.

A special case of saturation occurs when incompatible constraints are detected. Then, the whole abstract state is reduced to  $\perp$ , following the principles of reduced product [19]. As an example, when the abstract state contains the constraints  $S = [\alpha]$  and  $\mathbf{len}_S = 0$ , such a reduction is performed.

To summarize, the computation of  $\mathbf{guar}\mathbf{d}_{\Sigma}(C, \sigma^\#)$  involves the addition to  $\sigma^\#$  of a set of constraints that are derived from  $C$ . It is conservative in general. The termination of this computation follows from the fact that the added constraints only involve syntactic subcomponents of the elements of  $C$  and  $\sigma_s^\#$ .

*Example 1.* We consider the abstract state of Figure 3(b) and the constraint  $S_1 = S_l.[\delta]$  where  $S_1$  is a new symbolic sequence variable. First, the constraint is added to the abstract state. Second, compaction replaces the pattern  $S_l.[\delta]$  with  $S_1$  in all other constraints. Third, the numerical inequality  $\alpha \geq \mathbf{max}_{S_1}$  and the sortedness of  $S_1$  entails that  $S_1$  is sorted. Then,  $S_1 = \mathbf{sort}(S_1)$  implies that  $\delta$  is the maximum value of  $S_1$ . Moreover, the fact that  $S_l$  is a subsequence of  $S_1$  entails that  $\mathbf{min}_{S_1} \leq \mathbf{min}_{S_l}$  and  $\mathbf{max}_{S_l} \leq \mathbf{max}_{S_1}$ . Finally, since  $\delta \leq \mathbf{i}$ ,  $\mathbf{guar}\mathbf{d}_{\Sigma}$  also derives  $\mathbf{max}_{S_1} \leq \mathbf{i}$ . Finally,  $\mathbf{guar}\mathbf{d}_{\Sigma}$  produces:

$$S = S_1.S_r \wedge S_1 = S_l.[\delta] \wedge S = \mathbf{sort}(S) \wedge S_i = \mathbf{sort}(S_i), i \in \{l, r, 1\} \\ \wedge \mathbf{max}_{S_l} \leq \mathbf{max}_{S_1} = \delta \leq \mathbf{max}_{S_r} \wedge \mathbf{min}_{S_1} \leq \mathbf{max}_{S_l} \wedge \delta \leq \mathbf{i} \wedge \alpha_0, \alpha_1 \neq \mathbf{0x0}$$

*Example 2.* In this example, we show the detection of mutually cyclic constraints. We consider the abstract state  $S_1 = S_2.S' \wedge S_2 = S''.S_3$ , and the addition of  $S_3 = S_1.S'''$ . Inlining definition constraints for  $S_2$  and  $S_3$  would produce the cyclic

constraint  $S_1 = S'' . S_1 . S''' . S'$ . Thus, this also implies that  $S'$ ,  $S''$ ,  $S'''$  are empty and that  $S_1$ ,  $S_2$  and  $S_3$  are equal. After removal of the cycle,  $\text{guard}_{\Sigma}$  produces:

$$S' = S'' = S''' = [] \wedge S_1 = S_2 = S_3 \wedge S_1 = S_2 . S' \wedge S_2 = S'' . S_3$$

The operator  $\text{guard}_{\Sigma}$  is sound in the following sense:

**Theorem 1 (Soundness of  $\text{guard}_{\Sigma}$ ).** *For all abstract state  $\sigma^{\#}$  and constraint  $C$ , we have  $\{\sigma \in \gamma_{\Sigma}(\sigma^{\#}) \mid \sigma \models_s C\} \subseteq \gamma_{\Sigma}(\text{guard}_{\Sigma}(C, \sigma^{\#}))$ .*

*Verification of a sequence constraint.* Second, we define the *constraint verification operator*  $\text{verif}_{\Sigma} : \Sigma^{\#} \times \mathbb{C} \rightarrow \{\mathbf{false}, \mathbf{true}\}$  which inputs a constraint  $C$  and an abstract state  $\sigma^{\#}$  and returns **true** when it can prove that  $\sigma^{\#}$  entails  $C$ . It is conservative in the sense that it may return **false** even when the constraint is satisfied. The computation of  $\text{verif}_{\Sigma}(C, (\sigma_n^{\#}, \sigma_m^{\#}, \sigma_s^{\#}))$  proceeds as follows:

1. If  $\sigma_s^{\#}$  is  $\perp$ , it returns **true**.
2. For definition constraints  $S = E$ ,  $\text{verif}_{\Sigma}$  inlines the definitions of variables, and returns **true** when both sides rewrite into syntactically equal expressions. The absence of cyclic constraints ensures this exploration terminates.
3. Otherwise, it returns **false**.

For constraints of the form  $S = \mathbf{sort}(E)$ , the operator uses a specific rule (shown below) since variables inside the **sort** function may be arbitrarily reordered. Instead, we take advantage of the multi-set abstract domain to establish that  $S$  and  $E$  have the same contents.

$$\frac{S = \mathbf{sort}(S) \quad \text{multi}_S = \tau_{\text{mul}}(E)}{S = \mathbf{sort}(E)}$$

The operator  $\text{verif}_{\Sigma}$  is sound in the following sense:

**Theorem 2 (Soundness of  $\text{verif}_{\Sigma}$ ).** *For all abstract state  $\sigma^{\#}$  and constraint  $C$ , if  $\text{verif}_{\Sigma}(\sigma^{\#}, C) = \mathbf{true}$  then, we have  $\gamma_{\Sigma}(\sigma^{\#}) \subseteq \{\sigma \in \Sigma \mid \sigma \models_s C\}$ .*

### 3.3 Lattice operations

We now discuss join, widening and inclusion checking operations for loop analysis. We assume that  $\mathbb{D}_n^{\#}$  provides a conservative inclusion test operator  $\text{is\_le}_n$  (it inputs two elements of  $\sigma_n^{\#}$  and returns **true** only when it succeeds proving the first is included in the second), an over-approximate join operator  $\text{join}_n$  and a widening  $\text{widn}_n$ , and that  $\mathbb{D}_m^{\#}$  provides similar operators  $\text{is\_le}_m$ ,  $\text{join}_m$ , and  $\text{widn}_m$ , and we build similar operators for  $\Sigma^{\#}$ .

*Inclusion checking.* The inclusion test operator inputs two abstract states and returns a boolean. When it returns **true**, the concretization of the first abstract state is included into that of the second one. The inclusion checking algorithm is based on the constraint representation of abstract states and boils down to a repeated application of  $\text{verif}_{\Sigma}$ .

**Definition 3 (Inclusion checking operator).** The operator  $\text{is\_lc}_{\Sigma} : \Sigma^{\#} \times \Sigma^{\#} \rightarrow \{\text{true}, \text{false}\}$  is defined by:

$$\begin{aligned} \text{is\_lc}_{\Sigma}((\sigma_{n,0}^{\#}, \sigma_{m,0}^{\#}, \sigma_{s,0}^{\#}), (\sigma_{n,1}^{\#}, \sigma_{m,1}^{\#}, \wedge_i C_i)) \\ := \text{is\_lc}_n(\sigma_{n,0}^{\#}, \sigma_{n,1}^{\#}) \wedge \text{is\_lc}_m(\sigma_{m,0}^{\#}, \sigma_{m,1}^{\#}) \wedge (\wedge_i \text{verify}_{\Sigma}(C_i, \sigma_{s,0}^{\#})) \end{aligned}$$

The soundness of  $\text{is\_lc}_{\Sigma}$  follows from that of  $\text{is\_lc}_n$ ,  $\text{is\_lc}_m$ , and  $\text{verify}_{\Sigma}$ :

**Theorem 3.** The operator  $\text{is\_lc}_{\Sigma}$  is sound in the sense that, for all  $\sigma_0^{\#}, \sigma_1^{\#} \in \Sigma^{\#}$ , if  $\text{is\_lc}_{\Sigma}(\sigma_0^{\#}, \sigma_1^{\#}) = \text{true}$ , then  $\gamma_{\Sigma}(\sigma_0^{\#}) \subseteq \gamma_{\Sigma}(\sigma_1^{\#})$ .

*Upper bounds.* As usual, we define two over-approximate upper-bound operators, namely, a classical join operator  $\text{join}_{\Sigma} : \Sigma^{\#} \times \Sigma^{\#} \rightarrow \Sigma^{\#}$  and a widening  $\text{widen}_{\Sigma} : \Sigma^{\#} \times \Sigma^{\#} \rightarrow \Sigma^{\#}$  that ensures termination.

Essentially, the  $\text{join}_{\Sigma}$  operator proceeds component-wise (like  $\text{is\_lc}_{\Sigma}$  as defined in Definition 3) and essentially preserves sequence constraints that appear in both arguments. In the case of definition constraint, it first saturates the conjunctions of constraints, so as to maximize the possible sets of common constraints. The algorithm of  $\text{widen}_{\Sigma}$  is similar, except that it does not saturate its left argument for the sake of termination. This implies that  $\text{widen}_{\Sigma}$  always returns a conjunction of constraints that forms a subset of the constraints of its left argument.

Both operators are sound and furthermore,  $\text{widen}_{\Sigma}$  guarantees termination.

**Theorem 4 (Soundness of  $\text{join}_{\Sigma}$  and  $\text{widen}_{\Sigma}$ , termination of  $\text{widen}_{\Sigma}$ ).** For all abstract states  $\sigma_0^{\#}, \sigma_1^{\#}$ , we have:

$$\gamma_{\Sigma}(\sigma_0^{\#}) \cup \gamma_{\Sigma}(\sigma_1^{\#}) \subseteq \gamma_{\Sigma}(\text{join}_{\Sigma}(\sigma_0^{\#}, \sigma_1^{\#})) \quad \gamma_{\Sigma}(\sigma_0^{\#}) \cup \gamma_{\Sigma}(\sigma_1^{\#}) \subseteq \gamma_{\Sigma}(\text{widen}_{\Sigma}(\sigma_0^{\#}, \sigma_1^{\#})).$$

Moreover, the operator  $\text{widen}_{\Sigma}$  ensures termination, that is, for all sequence  $(\sigma_n^{\#})_{n \in \mathbb{N}}$  of abstract states the sequence  $((\sigma_n^{\#})'_{n \in \mathbb{N}})$  defined by  $(\sigma_0^{\#})'_0 = \sigma_0^{\#}$  and  $(\sigma_n^{\#})'_{n+1} = \text{widen}_{\Sigma}((\sigma_n^{\#})'_n, \sigma_{n+1}^{\#})$  is ultimately stationary.

*Example 3 (Join).* In this example, we consider the computation of the join of two abstract states taken from the analysis of the program of Figure 2. The analysis of the loop at line 7 involves the computation of the join of the three abstract states below. For concision, we omit inequality constraints involving extremal values of empty sequences.

$$\begin{aligned} \sigma_0^{\#} &::= \left\{ \begin{array}{l} S = S_0 \wedge S_1 = S_2 = [] \\ \wedge S = \text{sort}(S) \wedge S_i = \text{sort}(S_i), i \in \{0, 1, 2\} \end{array} \right. \\ \sigma_1^{\#} &::= \left\{ \begin{array}{l} S = S_0.S_2 \wedge S_1 = [] \wedge \max_{S_0} \leq \min_{S_2} \wedge i \leq \min_{S_2} \\ \wedge S = \text{sort}(S) \wedge S_i = \text{sort}(S_i), i \in \{0, 1, 2\} \end{array} \right. \\ \sigma_2^{\#} &::= \left\{ \begin{array}{l} S = S_1.S_0 \wedge S_2 = [] \wedge \max_{S_1} \leq \min_{S_0} \wedge \max_{S_1} \leq i \\ \wedge S = \text{sort}(S) \wedge S_i = \text{sort}(S_i), i \in \{0, 1, 2\} \end{array} \right. \end{aligned}$$

The most notable step is the saturation of the first argument, that injects constraint  $S = S_1.S_0.S_2$ , as a consequence of  $S = S_0$  and  $S_1 = S_2 = []$  in  $\sigma_0^{\#}$ ,  $S = S_0.S_2$  and  $S_1 = []$  in  $\sigma_1^{\#}$  and  $S = S_1.S_0$  and  $S_2 = []$  in  $\sigma_2^{\#}$ . After this, constraints that hold in only either argument are dropped, as, e.g., constraint  $S_1 = []$  in  $\sigma_1^{\#}$ . The result of the union corresponds to the abstract state in Figure 3(c).

```

struct list { struct list* next; int data; };

lsegs(α0, α1, S □) :=
| emp ∧ α0 = α1 ∧ S = □
| ∃ α', δ, S', α0.next ↦ α' * α0.data ↦ δ * lsegs(α', α1, S' □) ∧ α0 ≠ 0x0 ∧ S = [δ].S'

```

Fig. 6. A C list data-type and the inductive summarizing predicate describing list segments

## 4 Combination of sequence abstraction and shape analysis

In this section, we define a shape analysis with inductive predicates that infers invariants about both the layout of data-structures and the sequences of values they store. For the sake of simplicity, we consider only a singly-linked list predicate (Figure 6) throughout this section, although our analysis and its implementation are parameterized by user-defined inductive predicates [15,16]. The generalization to other structures will be discussed in Section 5.

### 4.1 Language and semantics

Although our implementation is based on the MemCAD analyzer [40] and targets the C language, our formalization only considers a restricted fragment. We let  $\mathbb{X}$  denote a finite set of program variables. We consider a basic imperative language, where commands are assignments, conditional statements, loops, and sequences of commands. Expressions are either l-values that evaluate to addresses, or r-values, that evaluate to scalars. An l-value  $l$  is either a program variable  $v \in \mathbb{X}$ , the access to an l-value field  $l.f$  (for concision, we let  $f$  denote both the field name and the corresponding memory offset), or the dereference  $*e$  of an expression  $e$ . An r-value  $e$  is either a constant  $n \in \mathbb{V}$ , or the reading of the memory cell defined by an l-value  $l$ , or the address  $\&l$  or an l-value  $l$ , or the application  $e_0 \oplus e_1$  of a binary operator to two sub-expressions. For simplicity, we assume here that operators are deterministic and cause no errors. The grammar is shown below:

$$l ::= v \mid l.f \mid *e \quad e ::= n \mid l \mid \&l \mid e \oplus e \quad c ::= l = e \mid \mathbf{if}(e)\{c\} \mid \mathbf{while}(e)\{c\} \mid c; c$$

We note  $\mathbb{A}$  for the set of addresses, which is a subset of the set of values  $\mathbb{V}$ . A memory state  $m$  is a partial function from addresses to values. We note  $\mathbb{M}$  for the set of memory states and let  $\emptyset$  denote the empty memory. Furthermore, we assume that each program variable  $x$  has a fixed address denoted by  $\underline{x} \in \mathbb{A}$ . Based on these definitions, we set up the program semantics as follows. First, we define the semantics of expressions by induction over their syntax. The semantics of an l-value  $l$  is a function  $\llbracket l \rrbracket_l : \mathbb{M} \rightarrow \mathbb{A}$  that maps a memory state  $m$  to the address  $l$  evaluates to in  $m$ . Similarly, the semantics  $\llbracket e \rrbracket_e : \mathbb{M} \rightarrow \mathbb{V}$  of an expression  $e$  maps a memory state to a value. Finally, the semantics  $\llbracket c \rrbracket : \mathcal{P}(\mathbb{M}) \rightarrow \mathcal{P}(\mathbb{M})$  of a command  $c$  maps any set of input memory states  $M$  to the set of all possible output memory states when starting from any  $m \in M$ . The definition of all three semantics is classical and shown in Figure 7, where  $f_{\oplus} : \mathbb{V}^2 \rightarrow \mathbb{V}$  denotes the semantics of operator  $\oplus$ .

$$\begin{aligned}
\llbracket \mathbf{x} \rrbracket_l(m) &:= \mathbf{x} & \llbracket n \rrbracket_e(m) &:= n \\
\llbracket l.\mathbf{f} \rrbracket_l(m) &:= \llbracket l \rrbracket_l(m) + \mathbf{f} & \llbracket l \rrbracket_e(m) &:= m(\llbracket l \rrbracket_l(m)) \\
\llbracket *e \rrbracket_l(m) &:= \llbracket e \rrbracket_e(m) & \llbracket \&l \rrbracket_e(m) &:= \llbracket l \rrbracket_l(m) \\
& & \llbracket e_0 \oplus e_1 \rrbracket_e(m) &:= f_{\oplus}(\llbracket e_0 \rrbracket_e(m), \llbracket e_1 \rrbracket_e(m)) \\
\llbracket l = e \rrbracket(M) &:= \{m \mid \llbracket l \rrbracket_l(m) \mapsto \llbracket e \rrbracket_e(m) \mid m \in M\} \\
\llbracket \mathbf{if}(e)\{c_0\} \rrbracket(M) &:= \llbracket c_0 \rrbracket(\{m \in M \mid \llbracket e \rrbracket_e(m) \neq 0\}) \cup \{m \in M \mid \llbracket e \rrbracket_e(m) = 0\} \\
\llbracket \mathbf{while}(e)\{c\} \rrbracket(M) &:= \{m \in \mathbf{lfp}F \mid \llbracket e \rrbracket_e(m) = 0\} \\
&\quad \text{where } F(M') = M \cup \llbracket c \rrbracket(m \in M' \mid \llbracket e \rrbracket_e(m) \neq 0) \\
\llbracket c_0; c_1 \rrbracket(M) &:= \llbracket c_1 \rrbracket \circ \llbracket c_0 \rrbracket(M)
\end{aligned}$$

**Fig. 7.** Semantics of programs

## 4.2 Combined memory and sequence abstraction

*Sequence aware shape abstraction.* We start with the definition of abstract memory predicates, following an approach similar to that of separation logic based shape analyses with inductive definitions [13,15], extended with sequence information. As explained early in the section, our formalization considers a single inductive predicate describing list segments, and parameterized with a symbolic sequence variable that stands for the sequence of the values contained in them (Figure 6). Considering only list segments has two advantages. First, complete lists can be expressed as list segments the last element of which has a “next” field equal to  $\mathbf{0x0}$ . Second, it simplifies reasoning over sequences as it avoids branching structures (considered in Section 5). Abstract states rely on scalar symbolic variables in  $\mathbb{X}_n$  to denote values and addresses and consist of separating conjunctions [52] of points-to predicates and of list segment predicates:

**Definition 4 (Abstract memory states).** *The set of abstract memory states  $\mathbb{M}^\sharp$  is described by the grammar below, where  $\alpha_0, \alpha_1 \in \mathbb{X}_n$  and  $S \in \mathbb{X}_s$ :*

$$m^\sharp ::= \mathbf{emp} \mid m^\sharp * m^\sharp \mid \alpha_0.\mathbf{f} \mapsto \alpha_1 \mid \mathbf{lseg}_s(\alpha_0, \alpha_1, S \square)$$

We note  $\mathbb{M}^\sharp$  for the set of abstract memory states.

As usual,  $\mathbf{emp}$  denotes the empty memory region and  $m_0^\sharp * m_1^\sharp$  denotes the disjoint union of memory regions described by  $m_0^\sharp$  (resp.,  $m_1^\sharp$ ). The abstract predicate  $\alpha_0.\mathbf{f} \mapsto \alpha_1$  denotes a single memory cell, the address of which is described by  $\alpha_0$  plus the offset of  $\mathbf{f}$  and the contents of which is described by  $\alpha_1$ . Finally,  $\mathbf{lseg}_s(\alpha_0, \alpha_1, S)$  stands for a (possibly empty) list segment that starts at an address described by  $\alpha_0$ , ending with a pointer to address  $\alpha_1$ , where each list element consists of two fields, namely, a pointer to the next element and a data field, and such that the sequence of the values of the data fields is described by sequence variable  $S$ . In logical terms, the predicate  $\mathbf{lseg}_s(\alpha_0, \alpha_1, S)$  is defined inductively as shown in Figure 6.

As the definition of  $\mathbf{lseg}_s$  in Figure 6 shows, the concretization of abstract memory states indirectly involves sequence variables (and also multi-set variables). Indeed, given an abstract memory state  $m^\sharp$  and a sequence variable  $S$  that appears in  $m^\sharp$ , the concretization of  $m^\sharp$  also constrains  $S$ ,  $\mathbf{len}_S$ , and  $\mathbf{multi}_S$ . To reflect

$$\begin{array}{c}
\frac{}{\emptyset, \sigma \models_{\mathbb{M}} \mathbf{emp}} \quad \frac{m = [\sigma_n(\alpha_0) + \mathbf{f} \mapsto \sigma_n(\alpha_1)]}{m, (\sigma_n, \sigma_m, \sigma_s) \models_{\mathbb{M}} \alpha_0.\mathbf{f} \mapsto \alpha_1} \quad \frac{\forall i, m_i, \sigma \models_{\mathbb{M}} m_i^\sharp}{m_0 \uplus m_1, \sigma \models_{\mathbb{M}} m_0^\sharp * m_1^\sharp} \\
\frac{\sigma_n(\alpha_0) = \sigma_n(\alpha_1) \quad \sigma_n, \sigma_s \models_s S = \square}{\emptyset, (\sigma_n, \sigma_m, \sigma_s) \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_0, \alpha_1, S)} \\
\frac{m, (\sigma_n, \sigma_m, \sigma_s) \models_{\mathbb{M}} \alpha_0.\mathbf{next} \mapsto \alpha_2 * \alpha_0.\mathbf{data} \mapsto \alpha_3 * \mathbf{lseg}_s(\alpha_2, \alpha_1, S_1 \square) \quad \sigma_n(\alpha_0) \neq 0 \quad \sigma_n, \sigma_s \models_s S = [\alpha_3].S_1 \quad S_1 \text{ fresh}}{m, (\sigma_n, \sigma_m, \sigma_s) \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_0, \alpha_1, S \square)}
\end{array}$$

**Fig. 8.** Concretization of abstract memory states

this, we let the concretization of an abstract memory  $m^\sharp$  return a set of tuples that comprise not only a memory state  $m$ , but also a valuation that maps each symbolic variable in  $m^\sharp$  to a value of the corresponding type (scalar, multi-set, or sequence). Such a valuation boils down to a triple  $(\sigma_n, \sigma_m, \sigma_s)$  (Definition 2). The definition of the concretization is based on a set of inductive derivation rules that follow the syntax of abstract memories and unfold the list segment predicates (Figure 8).

**Definition 5 (Concretization of abstract memory states).** *The concretization of abstract memory states  $\gamma_{\mathbb{M}}$  maps an abstract memory  $m^\sharp$  to a set of pairs  $(m, \sigma) \in \mathbb{M} \times \Sigma$  and is defined by:*

$$\gamma_{\mathbb{M}}(m^\sharp) = \{(m, \sigma) \mid (m, \sigma) \models_{\mathbb{M}} m^\sharp\}$$

As examples of abstract memory states, we refer the reader to the left conjuncts of the three abstract states shown in Figure 3.

*Combined abstract domain.* The analysis needs to reason accurately over sequence variables not only when they are bound in an inductive predicate, but also when these predicates are unfolded. Thus, it requires a product abstract domain based on the memory abstract domain fixed in Definition 4 and Definition 5 and on the sequence abstract domain introduced in Section 3. Moreover, like most shape analyses, it sometimes needs to make case splits due to the disjunctive nature of the inductive predicate  $\mathbf{lseg}_s$ . Thus, the combined abstraction is defined as follows:

**Definition 6 (Combined abstraction).** *The elements of the combined state abstract domain  $\mathbb{S}^\sharp$  are finite disjunctions of pairs of the form  $(m^\sharp, \sigma^\sharp) \in \mathbb{M}^\sharp \times \Sigma^\sharp$ . Furthermore, the concretization  $\gamma_{\mathbb{S}}$  maps an element  $s^\sharp$  of  $\mathbb{S}^\sharp$  into a set of memories  $m$  and is defined by:*

$$\gamma_{\mathbb{S}}((m^\sharp, \sigma^\sharp)) := \{m \mid \exists \sigma \in \gamma_{\Sigma}(\sigma^\sharp), (m, \sigma) \in \gamma_{\mathbb{M}}(m^\sharp)\} \quad \gamma_{\mathbb{S}}(\bigvee_i s_i^\sharp) := \bigcup_i \gamma_{\mathbb{S}}(s_i^\sharp)$$

*Concatenating segments.* Before we move to the analysis algorithms, we discuss a principle for logical reasoning over segments that many analysis operations rely on. Intuitively, a pair of consecutive segments may be merged into a single segment, that stores a sequence of elements that is the concatenation of the elements in the two initial segments. Reciprocally, it is possible to split a segment based on a partition of the sequence of its elements. The lemma below formalizes this.



**Lemma 1 (Concatenation (list predicates)).** *We assume  $\alpha_0, \alpha_1, \alpha_2$  distinct symbolic variables and let  $m_0^\# := \mathbf{lseg}_s(\alpha_0, \alpha_1, S_1 \square) * \mathbf{lseg}_s(\alpha_1, \alpha_2, S_2 \square)$ ,  $m_1^\# := \mathbf{lseg}_s(\alpha_0, \alpha_2, S \square)$ , and  $\sigma^\# := S = S_1.S_2$ . Then:*

- $\gamma_{\mathbb{S}}(m_0^\#, \sigma^\#) \subseteq \gamma_{\mathbb{S}}(m_1^\#, \sigma^\#)$ ;
- if  $(m, \sigma_1) \in \gamma_{\mathbb{S}}(m_1^\#, \sigma^\#)$ , then there exists  $\sigma_0$  such that  $(m, \sigma_0) \in \gamma_{\mathbb{S}}(m_0^\#, \sigma^\#)$  and, for all  $\beta \in \mathbb{V}$  such that  $\beta \neq \alpha_1$ ,  $\sigma_0(\beta) = \sigma_1(\beta)$ .

### 4.3 Computation of abstract post-conditions

Abstract post-conditions are computed by a pair of families of functions:

- given l-value  $l$  and expression  $e$ ,  $\mathbf{assign}_{\mathbb{S}, l=e} : \mathbb{S}^\# \rightarrow \mathbb{S}^\#$  computes an over-approximation for the assignment command  $l = e$ ;
- given expression  $e$ ,  $\mathbf{guard}_{\mathbb{S}, e} : \mathbb{S}^\# \rightarrow \mathbb{S}^\#$  computes an over-approximation for the effect of the condition expression  $e$ .

In the following paragraphs, we give the main steps of the algorithms to compute them. They both ensure the soundness conditions below, for all l-value  $l$ , expression  $e$ , and abstract state  $s^\# \in \mathbb{S}^\#$ :

$$\begin{aligned} \llbracket l = e \rrbracket(\gamma_{\mathbb{S}}(s^\#)) &\subseteq \gamma_{\mathbb{S}}(\mathbf{assign}_{\mathbb{S}, l=e}(s^\#)) \\ \{m \in \gamma_{\mathbb{S}}(s^\#) \mid \llbracket e \rrbracket_e(m) \neq 0\} &\subseteq \gamma_{\mathbb{S}}(\mathbf{guard}_{\mathbb{S}, e}(s^\#)) \end{aligned}$$

*Simple cases.* The computation of post-conditions for assignments and tests that involve only fully exposed cells is straightforward and follows classical shape analysis techniques [16]. For instance:

$$\begin{aligned} \mathbf{assign}_{\mathbb{S}, \underline{x}.f=y}(\underline{x}.f \mapsto \alpha_0 * \underline{y} \mapsto \alpha_1 * m^\#, (\sigma_n, \sigma_m, \sigma_s)) \\ = (\underline{x}.f \mapsto \alpha_1 * \underline{y} \mapsto \alpha_1 * m^\#, (\sigma_n, \sigma_m, \sigma_s)) \\ \mathbf{guard}_{\mathbb{S}, \underline{x}.f \neq 0 \times 0}(\underline{x}.f \mapsto \alpha_0 * m^\#, (\sigma_n, \sigma_m, \sigma_s)) \\ = (\underline{x}.f \mapsto \alpha_0 * m^\#, (\mathbf{guard}_n(\alpha_0 \neq 0, \sigma_n), \sigma_m, \sigma_s)) \end{aligned}$$

where  $\mathbf{guard}_n$  denotes a sound condition test for the numerical domain [21].

*Unfolding inductive predicates.* The more difficult cases in post-conditions arise when some of the memory cells that are affected by the statement are summarized as part of an inductive predicate as, e.g., in  $\mathbf{assign}_{\mathbb{S}, x=x.\mathbf{next}}(\underline{x} \mapsto \alpha_0 * \mathbf{lseg}_s(\alpha_0, \alpha_1, S \square))$ . In such cases, some inductive predicates need to be *unfolded*, before falling back to the simpler situation shown in the two aforementioned cases.

The unfolding operation is based on rewriting rules that follow directly from the inductive nature of  $\mathbf{lseg}_s$ . We note  $\rightsquigarrow$  the unfolding relation that rewrites an abstract state into another. Basic cases of  $\rightsquigarrow$  proceed as follows:

$$\begin{aligned} &(\mathbf{lseg}_s(\alpha_0, \alpha_1, S \square) * m^\#, (\sigma_n^\#, \sigma_m^\#, \sigma_s^\#)) \\ \rightsquigarrow &\left\{ \begin{array}{l} (m^\#, \mathbf{guard}_{\underline{x}}(S = \square), \mathbf{guard}_n(\alpha_0 = \alpha_1, \sigma_n^\#), \sigma_m^\#, \sigma_s^\#) \\ \vee \left( \alpha_0.\mathbf{next} \mapsto \alpha_2 * \alpha_0.\mathbf{data} \mapsto \alpha_3 * \mathbf{lseg}_s(\alpha_2, \alpha_1, S_1 \square) * m^\#, \right. \\ \left. \mathbf{guard}_{\underline{x}}(S = [\alpha_3].S_1, \mathbf{guard}_n(\alpha_0 \neq 0, \sigma_n^\#), \sigma_m^\#, \sigma_s^\#) \right) \end{array} \right\} \end{aligned}$$

where  $\alpha_2, \alpha_3, S_1$  are fresh. Unfolding is proved sound by the rules of Figure 8 in the sense that, for all  $s_0^\#, s_1^\# \in \mathbb{S}^\#$ , if  $s_0^\# \rightsquigarrow s_1^\#$ , then  $\gamma_{\mathbb{S}}(s_0^\#) \subseteq \gamma_{\mathbb{S}}(s_1^\#)$ .

The soundness of `assign`<sub>ℳ</sub> and `guard`<sub>ℳ</sub> follows from that of the unfolding relation, from that of the assignment and condition test of the underlying abstract domains, and from the (straightforward) handling of the unfolded cases.

We remark that the main difference compared to baseline shape analyses is that unfolding produces additional predicates about the sequence variables, which are added into the sequence domain. In turn, the addition of these constraints may yield increased precision due to internal reduction.

#### 4.4 Computation of lattice operations

The lattice operations required for the analysis of loops comprise the conservative inclusion test and the over-approximation of concrete upper bounds. Moreover, the former is used in the definition of the latter. Again, the algorithms to compute them are based on those of classical shape analyses. Thus, we emphasize the extensions that are required to infer sequence information and refer the reader to [16] for a full description of shape abstraction inclusion and widening algorithms.

*Inclusion checking.* The inclusion test function performs a proof search to try to establish inclusion. Although the rule system actually used is more complex, the inclusion proof system can be summarized down to three basic principles. First, when two abstract states have the same abstract memory component, proving inclusion boils down to checking the inclusion in  $\mathbb{S}^\#$ . Second, when the left-hand side contains several inductive predicate instances that can be summarized into one in the right-hand side, the analysis tries to concatenate them using Lemma 1. Third, when the right-hand side can be unfolded and the left-hand side is included into one of the unfolded disjuncts, then the inclusion holds for the initial pair. The rules below formalize these three principles.

$$\frac{\frac{\text{is\_le}_{\mathbb{S}}(\sigma_l^\#, \sigma_r^\#) = \text{true}}{(m^\#, \sigma_l^\#) \sqsubseteq (m^\#, \sigma_r^\#)}}{\text{verify}_{\mathbb{S}}(\sigma_l^\#, S = S_1.S_2) = \text{true}}}{(\text{lseg}_s(\alpha, \beta, S_1 \sqsupset) * \text{lseg}_s(\beta, \delta, S_2 \sqsupset) * m_l^\#, \sigma_l^\#) \sqsubseteq (\text{lseg}_s(\alpha, \delta, S \sqsupset), \sigma_r^\#)} \\ \frac{s_r^\# \rightsquigarrow \vee_i \overbrace{(m_i^\#, \text{guard}_{\mathbb{S}}(\sigma_i^\#, C_i))}^{s_i^\#}}{\exists j, \text{verify}_{\mathbb{S}}(C_j, \sigma_l^\#) = \text{true} \wedge (m_l^\#, \sigma_l^\#) \sqsubseteq s_j^\#}{(m_l^\#, \sigma_l^\#) \sqsubseteq s_r^\#}$$

The `is_le`<sub>ℳ</sub> function takes two abstract states and attempts to construct a proof tree that establishes inclusion based on these principles. The main specificities of the product with a sequence abstract domain are the requirement for `is_le`<sub>ℳ</sub> to track sequence concatenation constraints and the use of the inclusion checking function of the sequence abstract domain. The soundness of `is_le`<sub>ℳ</sub> follows from the soundness of the shape inclusion algorithm and of the underlying domains operations:

**Theorem 5 (Soundness of  $\text{is\_le}_S$ ).** *For all  $s_0^\#, s_1^\# \in \mathbb{S}^\#$ , if  $\text{is\_le}_S(s_0^\#, s_1^\#) = \text{true}$  then  $\gamma_S(s_0^\#) \subseteq \gamma_S(s_1^\#)$ .*

*Join and Widening.* The cases of join and widening are more subtle, since these operators may need to introduce  $\text{lseg}_s$  predicates together with fresh symbolic sequence variables, and to infer accurate relations over these new variables. Indeed, these algorithms are based on the following two principles:

- when the memory components of the two arguments are equal, we use it for the shape specific part of the result;
- when the memory components of the two arguments differ, they need to be *weakened* by replacing memory fragments with novel instances of  $\text{lseg}_s$ , with fresh symbolic sequence variables, and by checking inclusion holds using  $\text{is\_le}_S$ .

To illustrate the second case, we consider the over-approximation of the two abstract states below:

$$\begin{aligned} s_0^\# &:= (\alpha_0.\text{next} \mapsto \alpha_1 * \alpha_0.\text{data} \mapsto \alpha_3 * \text{lseg}_s(\alpha_1, \alpha_2, S \sqsupset), \sigma_0^\#) \\ s_1^\# &:= (\text{lseg}_s(\alpha_0, \alpha_1, S \sqsupset) * \alpha_1.\text{next} \mapsto \alpha_2 * \alpha_1.\text{data} \mapsto \alpha_3, \sigma_1^\#) \end{aligned}$$

Clearly, the memory part of both states may be weakened to the same abstract memory  $\text{lseg}_s(\alpha_0, \alpha_2, S'' \sqsupset)$  where  $S''$  is fresh. This gives the shape specific part of the result. However, in the case of  $s_0^\#$ , this weakening holds under the constraint  $S'' = [\alpha_3].S$ , whereas it holds under the constraint  $S'' = S'.[\alpha_3]$  in the case of  $s_1^\#$ . Therefore, the sequence abstract states should be updated according to these two constraints before calling the corresponding operator in the sequence domain, which produces:

$$\text{join}_\Sigma(\text{guard}_\Sigma(S'' = [\alpha_3].S, \sigma_0^\#), \text{guard}_\Sigma(S'' = S'.[\alpha_3], \sigma_1^\#)).$$

Note that this weakening also generates numerical and multi-set constraints. This constraint synthesis issue is carried out by an extension of the inclusion checking algorithm that keeps track of the fresh variables introduced by the widening and accumulates constraints over these.

**Theorem 6 (Soundness of  $\text{join}_S$ ,  $\text{widen}_S$  and its termination).** *The upper bound operator  $\text{join}_S, \text{widen}_S : \mathbb{S}^\# \times \mathbb{S}^\# \rightarrow \mathbb{S}^\#$  are sound in the sense that, for all  $s_0^\#, s_1^\# \in \mathbb{S}^\#$ , then*

$$\gamma_S(s_0^\#) \cup \gamma_S(s_1^\#) \subseteq \gamma_S(\text{join}_S(s_0^\#, s_1^\#)) \quad \gamma_S(s_0^\#) \cup \gamma_S(s_1^\#) \subseteq \gamma_S(\text{widen}_S(s_0^\#, s_1^\#))$$

*Moreover,  $\text{widen}_S$  also ensures the termination property [18].*

#### 4.5 Static analysis of a simple language

The analysis of a command  $c$  is a function  $\llbracket c \rrbracket^\# : \mathbb{S}^\# \rightarrow \mathbb{S}^\#$  that over-approximates  $\llbracket c \rrbracket$ . It is defined by induction over the syntax in Figure 9. Note that the convergence of the sequence of abstract iterates follows from the termination property of  $\text{widen}_S$ , and the analysis uses  $\text{is\_le}_S$  to detect stabilization. For conditional statements, we

$$\begin{aligned}
\llbracket l = e \rrbracket^\sharp(s^\sharp) &:= \text{assign}_{\mathcal{S}, l=e}(s^\sharp) \\
\llbracket \text{if}(e)\{c_0\} \rrbracket^\sharp(s^\sharp) &:= \text{join}_{\mathcal{S}} \left( \llbracket c_0 \rrbracket^\sharp(\text{guard}_{\mathcal{S}, e \neq 0}(s^\sharp)), \text{guard}_{\mathcal{S}, e=0}(s^\sharp) \right) \\
\llbracket \text{while}(e)\{c\} \rrbracket^\sharp(s^\sharp) &:= \text{guard}_{\mathcal{S}, e=0}(\lim_n s_n^\sharp) \\
&\text{where } s_0^\sharp := s^\sharp \text{ and } s_{n+1}^\sharp := \text{widen}_{\mathcal{S}}(s_n^\sharp, \llbracket c \rrbracket^\sharp(\text{guard}_{\mathcal{S}, e \neq 0}(s_n^\sharp))) \\
\llbracket c_0; c_1 \rrbracket^\sharp(s^\sharp) &:= \llbracket c_1 \rrbracket^\sharp \circ \llbracket c_0 \rrbracket^\sharp(s^\sharp)
\end{aligned}$$

**Fig. 9.** Abstract interpretation of a command

analyze the two branches separately after assuming the corresponding constraint, and we merge the two resulting states using  $\text{join}_{\mathcal{S}}$ . It is sound (the proof of soundness is classical [16] and proceeds by induction over the syntax):

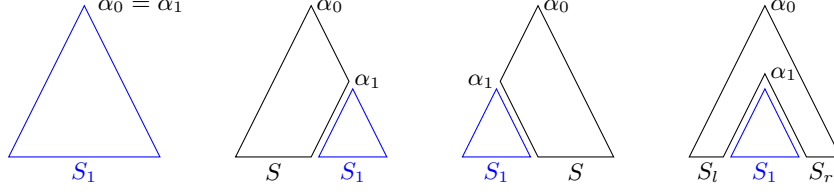
**Theorem 7 (Soundness).** *For all command  $c$ ,  $\llbracket c \rrbracket \circ \gamma_{\mathcal{S}} \subseteq \gamma_{\mathcal{S}} \circ \llbracket c \rrbracket^\sharp$ .*

## 5 Shape and sequence predicates for non-linear structures

This section discusses the general inductive predicates used by our analysis. While Section 4 only considered basic list predicates so as to introduce the analysis in a simpler setup, we now show our analysis handles  $m$ -ary trees (thus including lists when  $m = 1$ ), possibly with parent pointers. We require that the sequence arguments of inductive definitions denote (sub-)sets of elements stored in structure (we comment on this restriction in Remark 1). The following paragraphs show the specificities of sequence predicates for such data-structures, the derivation of segment predicates and how it affects analysis operations.

*Segment predicates and sequence information.* Segment predicates such as **lseg** play a very important role in the analysis, e.g., to analyze data-structure traversals, as in Section 2. Basic analysis operations split or merge inductive predicates that describe full structures and segments. As we remarked in Section 4, sequence information needs to be maintained when such steps are performed and Lemma 1 provides the method to do so for **lseg**. As observed in Section 2, the method derived from Lemma 1 will not work for non linear structures.

Indeed, let us consider the tree segment predicate **treeseq<sub>s</sub>** shown in Figure 4, which describes all the possible ways to decompose memory states that store a full tree at node  $\alpha_0$  and where  $\alpha_1$  is the address of one of its subtrees. Equivalently, the memory can be decomposed into a tree segment between  $\alpha_0$  and  $\alpha_1$  and a full tree at root  $\alpha_1$ . We note  $S_0$  (resp.,  $S_1$ ) the sequence of elements in the whole structure (resp., the subtree). Figure 10 depicts all possible configurations. In the first case, the subtree and the tree are equal, so the segment is empty and  $S_0 = S_1$ . In the second case, the subtree at  $\alpha_1$  is a leftmost subtree and  $S_0 = S_1.S_r$  for some  $S_r$ . The third case is symmetric. The fourth case is the most general and  $S_0 = S_l.S_1.S_r$  for some  $S_l, S_r$ . Therefore, the most general definition of the sequence(s) of elements in the segment (when the subtree in shown blue is excluded) is  $S_0 = S_l \sqsupset S_r$ , where  $\sqsupset$  is a placeholder that abstracts the sequence of the elements in the “missing” subtree, and where  $S_l, S_r$  may denote the empty sequence.



**Fig. 10.** Concatenation cases for tree segments and full tree predicates

Following this discussion, we now study concatenation of tree segment predicates. Let us assume two disjoint regions respectively abstracted by  $\mathbf{tree}_{\text{seg}}_s(\alpha, \alpha', S'_l \sqsupseteq S'_r)$  and by  $\mathbf{tree}_{\text{seg}}_s(\alpha', \alpha'', S''_l \sqsupseteq S''_r)$ . Then, the union of these two regions may be abstracted by  $\mathbf{tree}_{\text{seg}}_s(\alpha, \alpha'', S_l \sqsupseteq S_r)$  where  $S_l = S'_l.S'_l''$  and  $S_r = S''_r.S'_r$ . Note that the sequence expression attached to the latter segment is calculated as  $S_l.S'_l \sqsupseteq S''_r.S_r = (S_l \sqsupseteq S_r)[\sqsupseteq \leftarrow S'_l \sqsupseteq S'_r]$ . Similar reasoning may be carried out to concatenate a segment and a full tree predicate. Based on these observations, we propose a concatenation lemma for  $\mathbf{tree}_{\text{seg}}_s$ :

**Lemma 2 (Concatenation (tree case)).** *We assume symbolic variables  $\alpha, \alpha', \alpha''$  and sequence variables  $S, S', S_l, S_r, S'_l, S'_r, S''_l, S''_r$ .*

- Let  $m_0^\# := \mathbf{tree}_{\text{seg}}_s(\alpha, \alpha', S'_l \sqsupseteq S'_r) * \mathbf{tree}_s(\alpha', S')$ ,  $m_1^\# := \mathbf{tree}_s(\alpha, S)$   $\sigma^\# := S = S'_l.S'.S'_r$ . Then,  $\gamma_{\mathbb{S}}(m_0^\#, \sigma^\#) \subseteq \gamma_{\mathbb{S}}(m_1^\#, \sigma^\#)$ .
- Let  $m_0^\# := \mathbf{tree}_{\text{seg}}_s(\alpha, \alpha', S'_l \sqsupseteq S'_r) * \mathbf{tree}_{\text{seg}}_s(\alpha', \alpha'', S''_l \sqsupseteq S''_r)$ ,  $m_1^\# := \mathbf{tree}_{\text{seg}}_s(\alpha, \alpha'', S_l \sqsupseteq S_r)$   $\sigma^\# := S_l = S'_l.S'_l'' \wedge S_r = S''_r.S'_r$ . Then,  $\gamma_{\mathbb{S}}(m_0^\#, \sigma^\#) \subseteq \gamma_{\mathbb{S}}(m_1^\#, \sigma^\#)$ .

*Derivation of segment predicates from full predicates.* While inductive predicates (e.g., the definition of lists or trees) are user-supplied, our analysis automatically derives the corresponding segment predicates. Indeed, given a full predicate (like  $\mathbf{tree}_s$ ) for an  $m$ -ary form of tree (including lists), the segment predicate is obtained by the sequence of steps below:

- each sequence argument  $S_i$  is replaced by a marked sequence  $S_i \sqsupseteq S'_i$ ,
- a rule describes empty segments; it abstracts an empty memory region, constrains its extremal points to be equal and its sequence contents to be empty;
- for each inductive rule that contains recursive calls to the inductive predicate, and for each such call  $c$ , the segment predicate should include a rule replacing  $c$  with a segment instance; moreover, in each such segment rule, the linearity of the sequence concatenations should be reflected by sequence constraints.

As an example, we illustrate this in the case of  $\mathbf{tree}_s$ :

*Example 4 (Tree segments).* The definition of  $\mathbf{tree}_s$  is shown in Figure 1. As it has one sequence parameter, the corresponding segment predicate has two, that we note  $S_0$  and  $S_1$  and writes down  $\mathbf{tree}_{\text{seg}}_s(\alpha, \alpha', S_0 \sqsupseteq S_1)$ . We now detail the derivation of the  $\mathbf{tree}_{\text{seg}}_s$  predicate shown in Figure 4. As stated above,  $\mathbf{tree}_{\text{seg}}_s$  includes a rule for empty segments (the first one in Figure 4), which corresponds

to an empty region, two equal pointers and two empty sequences. The first rule of  $\mathbf{tree}_s$  corresponds to the empty tree; it has no recursive call and cannot appear in segments. The second rule of  $\mathbf{tree}_s$  has two recursive calls (for the left and right subtrees), thus, it gives rise to two rules in  $\mathbf{tree}_{\mathbf{seg}}_s$ , that stand for cases where the segment is in the left (resp., right) subtree. Finally, we consider the constraints over sequences in the last rule (right subtree). Given the notation in Figure 4, the sequence of values in the whole tree is the argument  $S \sqsupset S'$  of  $\mathbf{tree}_{\mathbf{seg}}_s$  which is equal to  $S_l.[v].S_r \sqsupset S'_r$  in the last rule. This equality entails the constraints  $S = S_l.[v].S_r$  and  $S' = S'_r$  which thus appear in the last rule of  $\mathbf{tree}_{\mathbf{seg}}_s$ .

*Remark 1 (Limitation of sequence arguments).* We observe that the inference of the sequence constraints by linearity as shown in Example 4 can only be achieved since the sequence constraints in  $\mathbf{tree}_s$  specify that its segment argument collects a set of elements found at some fields in the structure. As an example, the analysis would not support an alternative definition of  $\mathbf{tree}_s$  where the inductive rule would have the sequence constraint  $S = \mathbf{sort}(S_l.[v].S_r)$ , as it does not allow the derivation of precise constraints over sub-sequences for segments. We note that this limitation does not prevent capturing precisely binary search trees in the product abstract domain of Definition 5 with element  $\mathbf{tree}_s(\alpha, S) \wedge S = \mathbf{sort}(S)$ ; instead, it only requires the shape predicate be written in a certain way.

*Analysis.* The analysis requires users to supply inductive predicates for full structures as well as target pre- and post-conditions. Segment predicates are inferred automatically as shown in the previous paragraph, as well as the appropriate concatenation lemma. Finally, the analysis operators are similar to those shown in Section 4, except that they use the concatenation property inferred from the definition of the full structure inductive predicate. For instance, when using the tree inductive predicate of Figure 1, the analysis infers the segment of Figure 4 and the concatenation lemma 2. The analysis satisfies the soundness property of Theorem 7. To conclude the section, we discuss a couple of steps of the computation of widening in the analysis of the program in Figure 2.

*Example 5 (Inclusion checking).* We consider the following abstract states:

- $s_0^\sharp = (\alpha.l \mapsto \alpha_0 * \alpha.d \mapsto \alpha_1 * \alpha.r \mapsto \alpha_2 * \mathbf{tree}_s(\alpha_2, S_r), \sigma_0^\sharp)$ ;
- $s_1^\sharp = (\mathbf{tree}_{\mathbf{seg}}_s(\alpha, \alpha_0, S \sqsupset S'), \sigma_1^\sharp)$ .

Both  $s_0^\sharp$  and  $s_1^\sharp$  appear during the widening at the first iteration. We study the evaluation of the inclusion test  $\mathbf{is\_lc}_S(s_0^\sharp, s_1^\sharp)$ . We first remark the following unfoldings (where  $\alpha_3, \alpha_4, \alpha_5, S_l$ , and  $S'_l$  are fresh) yield a similar abstract memory, up to existentially quantified symbolic variable names:

$$\begin{aligned} s_1^\sharp &\rightsquigarrow (\alpha.l \mapsto \alpha_3 * \alpha.d \mapsto \alpha_4 * \alpha.r \mapsto \alpha_5 * \mathbf{tree}_{\mathbf{seg}}_s(\alpha_3, \alpha_0, S_l \sqsupset S'_l) \\ &\quad * \mathbf{tree}_s(\alpha_5, S_r), \sigma_0^\sharp) \wedge S = S_l \wedge S' = S'_l.[\alpha_4].S_r \\ &\rightsquigarrow (\alpha.l \mapsto \alpha_3 * \alpha.d \mapsto \alpha_4 * \alpha.r \mapsto \alpha_5 * \mathbf{emp} * \mathbf{tree}_s(\alpha_5, S_r), \\ &\quad \sigma_0^\sharp) \wedge S = S_l \wedge S' = S'_l.[\alpha_4].S_r \wedge S_l = [] \wedge S'_l = [] \wedge \alpha_0 = \alpha_3 \end{aligned}$$

By the definition of  $\mathbf{is\_lc}_S$  in Section 4.4,  $\mathbf{is\_lc}_S(s_0^\sharp, s_1^\sharp)$  returns true if and only if  $\mathbf{is\_lc}_\Sigma(\sigma_0^\sharp, \sigma_1^\sharp)$ ,  $\mathbf{verify}_\Sigma(\sigma_0^\sharp, S = [])$  and  $\mathbf{verify}_\Sigma(\sigma_0^\sharp, S' = [\alpha_1].S_r)$  all return true.

$$\begin{array}{c}
\underbrace{\left( \begin{array}{l} \&t \mapsto \alpha_0 \\ * \&c \mapsto \alpha_0 \\ * \mathbf{emp} \\ * \mathbf{tree}_s(\alpha_0, S) \end{array} \right)}_{m_0^\sharp} \quad \underbrace{\left( \begin{array}{l} \&t \mapsto \alpha_0 * \&c \mapsto \alpha_1 \\ * \alpha_0.l \mapsto \alpha_1 * \mathbf{tree}_s(\alpha_1, S_l) \\ * \alpha_0.d \mapsto \alpha_2 \\ * \alpha_0.r \mapsto \alpha_3 * \mathbf{tree}_s(\alpha_3, S_r) \end{array} \right)}_{m_1^\sharp} \quad \underbrace{\left( \begin{array}{l} \&t \mapsto \alpha \\ * \&c \mapsto \alpha' \\ * \mathbf{tree}_{\text{seg}_s}(\alpha, \alpha', S_1 \sqcap S_2) \\ * \mathbf{tree}_s(\alpha', S_0) \end{array} \right)}_{m_f^\sharp} \\
\begin{array}{|c|c|c|c|c|c|} \hline m_f^\sharp & \alpha & \alpha' & S_0 & S_1 & S_2 \\ \hline m_0^\sharp & \alpha_0 & \alpha_0 & S & \square & \square \\ \hline m_1^\sharp & \alpha_0 & \alpha_1 & S_l & \square & [\alpha_2].S_r \\ \hline \end{array}
\end{array}$$

**Fig. 11.** Shape union between states from Figures 3(a) and 3(b) (Greek letters denote existentially quantified symbolic variables; identical colors denote similar regions).

*Example 6 (Widening).* We now study the computation of the first widening in the analysis of the program shown in Section 2. For brevity, we only consider the second disjunct after the condition. The arguments of widening of abstract memory states and the result are shown in Figure 11. As mentioned in Section 4.4, the widening operator seeks for regions that can be described in a similar manner in the both of its arguments, possibly after weakening them. Matching colors in Figure 11 highlight pairings of similar regions. Recall that all symbolic variables  $(\alpha, \alpha_0, \dots)$  are existentially quantified within a same state. We observe the terms in blue, green and purple are pairwise equal and require no weakening. The areas in red though are not equal. For clarity, we add an **emp** term in  $m_0^\sharp$ . As observed in Example 5, the matching terms in  $m_1^\sharp$  can be weakened into  $\mathbf{tree}_{\text{seg}_s}(\alpha_0, \alpha_1, S_1 \sqcap S_2)$ , provided  $S_1 = \square$  and  $S_2 = [\delta].S_r$ . The same holds for **emp** in  $m_0^\sharp$ . The table in the bottom of Figure 11 summarizes the correspondence between existentially quantified symbolic variables that realizes the association of regions.

The above paragraph describes the computation of the abstract memory state shown in Figure 3(c). The computation of the sequence abstract state of Figure 3(c) proceeds by application of  $\mathbf{widen}_\Sigma$ .

## 6 Implementation and evaluation

In this section, we report on the implementation and evaluation of the product shape and sequence analysis. We consider the following research questions:

- **(RQ1)** Is the combined analysis of Section 4 and Section 5 precise enough to prove functional properties on programs implementing classical algorithms over dynamic data-structures (like lists, sorted lists, and binary search trees), and does it help a baseline analysis verify structural invariants are preserved?
- **(RQ2)** Can this analysis successfully verify real-world C libraries?
- **(RQ3)** How significant is the overhead of the combined analysis compared to the baseline?

*Implementation.* We have implemented the sequence abstract domain and the product with the shape abstraction of the MemCAD static analyzer [40,1]. The

Example	without seq			with seq parameters					PrSafe + Fc verified
	time	#iter	PrSafe verified	time				#iter	
	all			all	num	seq	shape		
Singly linked list									
Push	4.0		✓	4.8	0.5	0.5	0.9		✓
Pop	5.1		✓	5.4	0.9	1.4	0.8		✓
Pop (empty)	4.9		✓	4.7	0.8	0.5	1.4		✓
concat	6.5	2	✓	15.7	3.4	3.3	2.7	2	✓
deep copy	12.1	2	✓	20.4	3.7	2.9	5.5	2	✓
length	9.5	3	✓	45.0	22.5	5.0	8.1	3	✓
insert at position	19.0	3	✓	101.9	61.3	7.9	12.2	3	✓
remove at position	17.2	3	✓	92.5	55.5	6.5	12.5	3	✓
inserting in a sorted list	13.5	3	✓	82.5	39.0	10.0	9.2	3	✓
minimum	11.8	3	✓	92.3	42.4	11.1	16.8	3	✓
maximum	11.8	3	✓	93.2	42.9	11.2	17.0	3	✓
insertion sort	24.6	2, 2	✓	714.6	328.6	90.0	126.3	4, 3	✓
bubble sort	40.6	2;2,3	✓(†)	776.3	399.5	89.2	141.5	3;3,3	✓(†)
merge	36.8	4	✓	352.2	180.9	41.0	54.9	4	✓
Binary trees									
Delete leftmost	11.2	3	✓	80.5	38.2	9.4	12.0	3	✓
Delete rightmost	11.5	2	✓	58.1	27.5	6.8	7.6	2	✓
Binary search trees									
Insertion	25.2	2	✓	150.4	58.0	17.2	15.5	2	✓
Delete max	22.9	2	✗	141.2	68.6	15.2	17.2	2	✓
Delete min	22.0	3	✗	177.9	87.9	19.2	22.8	3	✓
Search (present)	26.6	2	✓	107.2	48.6	15.7	14.4	2	✓
Search (absent)	24.0	3	✓	76.7	29.4	11.4	11.7	3	✓
BST to list (heap sort)	23.8	3	✓	76.5	29.2	11.4	11.7	3	✓
list to BST (heap sort)	34.2	2,2	✓	408.0	188.0	56.5	68.4	3,2	✓

**Table 1.** Experimental results on custom examples (Time in milliseconds averaged over 100 runs. For loop iterations, disjoint loops are separated by a semicolon, nested loops by a comma, and the first number corresponds to the outer loop. For inner loops, we take the maximum number of iterations needed to stabilize it.)

analysis inputs C programs and user-supplied inductive predicates describing data-structures together with pre- and post-conditions and attempts to verify them, as well as absence of runtime errors. We set convex polyhedra [21] implemented in the Apron library [36] as numerical abstraction and an extension of [41] as multi-set abstraction.

*Experiments.* We consider two sets of experiments. The first one (Table 1) consists of custom implementations of classical algorithms over lists, sorted lists, and binary search trees and includes sorting, insertion and deletion algorithms. The second (Table 2) collects list data-structure implementations taken from the Linux [56] and FreeRTOS [34] operating systems as well as the Generic data-structure library (GDSDL) [24], which all involve specificities like back pointers or sentinel nodes. For each data-structure, we provide an inductive definition written in the DSL of MemCAD. This amounts to a single definition a few lines long for each series of tests. For each test, we also specify the pre- and post-condition of procedures. When a procedure may behave differently depending on the shape of their input, we provide two pre-/post-condition pairs. This occurs for the “Pop” function, which does nothing when applied to the empty list. Two target properties are studied:

- **PrSafe:** absence of memory errors and structural preservation (with respect to list or tree invariants but without checking anything about their contents);



Example	without seq			with seq parameters					
	time	#iter	PrSafe verified	time				#iter	PrSafe + Fc verified
	all			all	num	seq	shape		
Linux lists									
Init	1.1		✓	2.6	0.2	0.3	1.1		✓
Input	13.6		✓	21.4	2.7	2.4	8.2		✓
Output	22.7		✓	31.5	4.8	4.8	10.5		✓
Output (empty)	33.8		✓	9.3	1.4	1.0	2.5		✓
FreeRTOS lists									
vListInit	4.3		✓	6.1	1.3	0.4	0.6		✓
vListInsertEnd	23.8		✓	40.3	10.8	1.8	5.3		✓
vListInsert	87.4	4	✓	370.5	202.4	27.2	37.9	4	✓
vListRemove	47.5		✓	163.4	82.6	9.2	20.0		✓
GDSDL (lists)									
Flush	24.3	2	✓	59.4	18.4	5.4	16.1	2	✓
Free	35.3	2	✓(†)	79.9	25.1	7.4	24.0	2	✓(†)
Remove head (empty)	34.1		✓	111.9	50.9	6.5	25.4		✓
Remove head (non-empty)	34.0		✓	16.3	5.7	1.1	3.7		✓
Remove tail (empty)	49.5		✓	284.8	165.0	13.6	39.3		✓
Remove tail (non-empty)	49.5		✓	16.2	5.7	1.1	3.6		✓
Search max	69.7	5	✓	708.4	429.7	43.1	145.7	5	PrSafe Fc
Search min	69.4	5	✓	634.0	380.3	35.4	131.2	5	PrSafe Fc
Search by position	104.5	3;2	✗(†)	1182.8	796.3	40.7	108.2	3;3	✓(†)

**Table 2.** Experimental results on libraries programs (Time in milliseconds averaged over 100 runs. For loop iterations, disjoint loops are separated by a semicolon, nested loops by a comma, and the first number corresponds to the outer loop. For inner loops, we take the maximum number of iterations needed to stabilize it.)

- **Fc**: partial functional correctness (including sortedness and the preservation of the elements stored in data-structures).

We ran the experiments on a machine with an i7-8700 processor with 32 Gb of RAM running Ubuntu 18.04. For each test case, we run the analysis *without* and then *with* sequence abstraction to compare runtimes and check if the analyses prove the expected property. When using the analysis without sequence abstraction, only **PrSafe** is considered (this abstraction cannot express **Fc**), whereas the analysis of sequences attempts to discharge both **PrSafe** and **Fc**. Table 1 displays raw results for the first series of tests. Table 2 shows the results of the main tests in the second series of tests (for brevity, all results are given in appendix in Table 3).

*Verification of complex properties.* As shown in Table 1, the analysis with sequences fully verifies both memory safety and functional correctness (**PrSafe** and **Fc**) for all target codes including three different list sorting programs, operations on binary search trees as well as heap sort (elements of a list are all inserted in an empty binary search tree and collected in a left to right order back into a list). These examples all require the inference of fairly involved invariants. The analysis without sequences can only verify **PrSafe**, yet it fails to do so in several examples, where the use of sequences actually also lets the analysis verify **PrSafe** (in addition to **Fc**). This result is somewhat surprising, as we would not expect sequence information be required to establish basic safety. One caveat is that one example (bubble sort) required the manual insertion of a directive to MemCAD to delay folding. We conjecture the shape folding operator could be improved to avoid this. All other

analyses are fully automatic. We conclude the product with sequences not only allows to prove **Fc** even in challenging cases, but may also help with **PrSafe**.

*Verification of real-world libraries.* We now consider Table 2. These examples involve lists with invariants that are considerably more sophisticated than **lseg<sub>s</sub>**, as they are all doubly-linked lists with headers. While GDSL lists contain a pointer to stored value blocks, both Linux and FreeRTOS lists are intrusive lists in the sense of the Linux kernel terminology: the **C** struct containing the **next** and **prev** fields is a substructure of the list node, which implies structure accesses require more complex pointer operations. FreeRTOS lists explicitly store a pointer from substructures to owners, whereas Linux lists rely on pointer arithmetic to access containing blocks. Finally, both FreeRTOS and GDSL lists have a header that stores the number of elements in the lists. FreeRTOS list nodes store a pointer to this header. The analysis with sequences proves both **PrSafe** and **Fc** for all Linux and FreeRTOS primitives. It was also able to fully verify almost all the GDSL list library, although two cases required a manual directive to prevent aggressive folding both with and without sequences (as for bubblesort in Table 1) (they are marked with (†) in the tables). Only two functions for the extraction of minimal/maximal values could not be fully verified with respect to **Fc** (note that **PrSafe** still gets proved): in these codes, the memory widening is too aggressive and folds the node storing the function results, which prevents proving that the returned value is indeed the extremal value in the sequence. The examples only shown in Table 3 are all verified. We conclude the analysis handles real-world programs.

*Overhead.* We now compare performance between the analyses with/without sequences in Tables 1 and 2. While the overhead is modest for the smaller programs, it becomes higher for the more challenging cases, up to roughly 10x-20x. While significant, this cost should be considered in comparison to the much stronger properties proved (i.e., not only **PrSafe** but also partial correctness **Fc** in addition to **PrSafe**). We found two reasons for this increase. First, as shown in the tables, most of the increase is accounted for by the numerical abstract domain partly due to the larger number of symbolic variables that stand for sequence bounds. We believe this overhead could be much reduced with a finer-grained numerical domain packing [7,54]. By contrast, the time spent in memory and sequence domains remains reasonable. Second, the analysis with sequences requires greater numbers of abstract iterates to stabilize loop iterates, as shown in the tables, which explains an important slowdown. This is to be expected due to the more complex value constraints (including polyhedra) used in the analysis with sequences.

## 7 Related works

In this section, we discuss previous works on the abstractions of sequences stored in data-structures.

*Linear and contiguous structures (arrays and strings).* Several previous works have tried to tie properties of container data-structures with properties of their contents. In particular, [29,30] have extended array abstractions with basic contents properties. More recently [31] introduced array segmentations and [20] made the computation of the array segmentations dynamic during the analysis. The latter two can express that an array is sorted and verify that a function produces sorted arrays. However, they do so with specific predicates rather than an abstraction for sequences. Thus, they cannot express that the set of elements in an array is preserved, which is required to prove a sorting function correct. By contrast, our sequence abstraction handles both sortedness and contents preservation.

Strings and buffers also motivated many research works, as operations on them may incur a security risk. In particular, improper handling of zero terminated strings make opens the door to buffer overrun attacks. Therefore, works such as CSSV [26] abstract the presence or absence of zeroes in strings and their positions in order to verify buffer operations. Besides zeroes, these works do not keep any contents' information.

As noted earlier, several recent works applied concepts such as regular expressions and automata in order to build string abstract domains, that convey precise contents information [47,3,49]. These works are typically aimed at inferring precise information on strings that denote pieces of programs meant to be computed and evaluated at runtime as in the case of JavaScript's `eval` construction. Automata and regular expressions are most adequate for such target properties. More recently, [4] extended these works with length and element position constraints. These abstractions are not aimed at numerical sequences, and fail to express sortedness. By contrast, our sequence abstraction relies on length, extremal elements and sortedness constraints and fails to express regular expressions-based properties as these would not be useful for our intended application. An interesting area of future work would be to build a reduced product of sequence abstractions to combine the expressiveness of these works and of ours.

*Shape analyses for dynamic data-structures.* Many abstractions for dynamic data-structures have been proposed. Sagiv *et al.* introduced a shape analysis based on three value logic in [53], that was later extended to handle more complex data-structures such as tree [44]. The seminal work by Reynolds [52], introduced separation logic, that many analyses including ours rely upon. Separation logic has been used in order to reason over not only sequential programs [13] but also concurrent programs [50,59] and to prove properties like linearizability of concurrent data structures [58]. It serves as a basis for structure abstraction in several static analyzers like Smallfoot [13], Facebook Infer [14] (which also performs bi-abduction to synthesize pre- and post-condition pairs), Forester [32] (which uses automata to represent abstract states), and MemCAD [40] (which features a modularized abstract domain). Bi-abduction methods have also been extended to infer inductive predicates on a per-function basis [39] or to infer pre- and post-conditions for programs manipulating lists and using bit-level memory accesses and pointer arithmetic [33]. All the shape abstractions mentioned so far can only keep track of very limited contents properties.

Indeed, inferring precise information about the contents of dynamic data-structures is notoriously difficult, since the memory abstraction layout changes depending on the program point which makes abstraction complex. A first approach to this issue consists in splitting the analysis in two phases, where the first analysis infers only structural invariants and translates the initial program into a purely numerical program, that is taken as input by the second analysis, that discovers numerical invariants. This technique has been applied by [45,28] in order to infer complexity bounds and verify termination of programs based on information on the size of the data-structures. A second approach [15] consists of a reduced product between a memory abstract domain and a numerical abstract domain. While harder to implement, it ensures information can be communicated in both directions between the memory and the value abstract domains, whereas the staged analysis approach only lets the value abstract domain benefit from memory layout information. More recently, [41] combines shape and set abstractions with a reduced product which allows verifying programs on graphs. As it only considers set constraints, it does not capture any order information.

The tools CINV [8] and CELIA [9] (extended with interprocedural analysis support in [12]) are the most closely related to our approach. These static analyzers handle list manipulating programs and are parameterized by an abstract domain called a *data-word domain* to reason on the structure and contents of lists by attaching size or set constraints to them, or constraints quantified over the position of elements, which allows expressing sortedness. Although the heap abstraction does not make explicit use of separation logic the list abstraction follows a similar structure. A first important difference with our work is that CINV and CELIA only handle singly-linked lists, whereas our analysis supports a large range of inductive definitions included doubly linked-lists, trees, binary search trees with and without parent pointers. Indeed, our approach integrates sequence reasoning into a shape analysis that can be parameterized by a wide variety of inductive predicates. This more general scope requires extensions to the analysis algorithms, such as the automatic inference of concatenation lemmas (Lemma 1 and 2) and the use of abstract operators based on them. A second difference comes from the sequence domain and the interaction with it. The data-word domain to handle sortedness relies on a decidable fragment of first order array theory based on constraints of the form  $\forall \mathbf{y}, P(\mathbf{y}) \Rightarrow U(\mathbf{y}, Q_1, \dots)$ , where the guard constraint  $P(\mathbf{y})$  belongs to a predefined, user-provided set of *guard-patterns* constraining the index variables  $y_j$ , and  $U$  is a conjunction of linear constraints on  $y_j$  and  $Q_i[y_j]$ . This domain does not manipulate symbolic sequence expressions but rather follows a structural approach. For example, the concatenation constraint  $S = S_1.S_2$  is expressed as  $\forall y_1, y_2, y_1 < \text{len}_{S_1} \wedge y_2 < \text{len}_{S_2} \Rightarrow S_1[y_1] = S[y_1] \wedge S_2[y_2] = S[y_2 + \text{len}_{S_1}]$ . Therefore, it requires the user to specify prior to the analysis the appropriate guard pattern. Our sequence abstraction requires no such parameterization.

*Provers for memory and contents properties.* Separation logic has also been used as foundation for verification tools based on entailment checking procedures, some of which also consider contents properties. Songbird [55] uses a sequent-based approach to attempt deciding implication in a fragment of separation logic enriched

with pure predicates. The procedure presented in [35] relies on tree automata to decide implications that involve inductive predicates. CSL [11] and SLAD [10] decide entailment on a logic for singly-linked lists and the data stored in them. It handles order constraints on linear structures like lists and arrays. More recently, [23] used bi-abduction to reason about ordered data by explicitly storing bounds on elements in the inductive predicate. This work only considers full structure predicates and does not handle segment predicates. All these tools can be used to discharge implication proof obligations and can be used in verification tools where invariants are either manually written or inferred by some other means.

Additionally, separation logic is also heavily used in approaches based on proof assistants [17,37]. In that case, contents properties are naturally expressed in the proof assistant language.

*Solvers for sequence properties.* Finally, we remark that our language of sequence constraints based on concatenation of atoms has some similarity with the string logic that can be found in some decision procedures. Though the logic of word equations with at least two atoms is known to be undecidable [51], its quantifier free fragment has a PSPACE complete decision procedure [46]. Following the work of [2] that classifies the field of string constraints solving in three main branches, the automata based approach, using finite state automata to represent the set of constraints [42], the word based approach, that decomposes constraints using algebraic results such as Levi’s lemma [6], and the unfolding based approach, which expresses each string variable as a bounded sequence of variables such as bit vectors [38], our abstraction can be categorized as mostly word-based. To the best of our knowledge, no SMT solver is able to reason on the sortedness of word expressions. We refer the reader to [2] for a comprehensive survey on string constraint solving. By comparison with these works, we provide an abstract domain interface on top of the sequence operation, which allows its use in a static analysis tool, following an instance of reduced product [19].

## 8 Conclusion and future works

In this paper, we presented a novel sequence abstract domain that relies on existing numerical and set abstractions, and extended a shape analysis with sequence reasoning. We demonstrated that the resulting analysis can be used in order to verify not only memory safety or structural preservation but also far more advanced correctness properties on a wide variety of inductive structures including various kinds of lists and trees. In particular, it could prove the functional correctness of several list sorting programs and of operations over binary search trees.

A combination of our analysis with a termination analysis [57,28] could verify not only partial correctness but also full correctness, which would be a first interesting direction for future works. Defining a reduced product over abstract domains for sequences would be also a useful research direction, as it would allow to strengthen the expressiveness of the analysis. Last, the evaluation also shows that performance of the combined analysis could be improved with the use of a more efficient dynamic packing [54] for relational constraints.

**Acknowledgment** The authors want to thank Thierry Martinez for his work on libraries used by the MemCAD analyzer. This work was supported by the VeriAMOS ANR-18-CE25-0010 French ANR project.

## References

1. Artifact for "A Product of Shape and Sequence Abstractions". Zenodo (Jul 2023), <https://doi.org/10.5281/zenodo.8186871>
2. Amadini, R.: A survey on string constraint solving. *ACM Computing Survey* (2021)
3. Arceri, V., Mastroeni, I.: An automata-based abstract semantics for string manipulation languages. In: *VPT@Programming* (2019)
4. Arceri, V., Olliaro, M., Cortesi, A., Ferrara, P.: Relational string abstract domains. In: *VMCAI* (2022)
5. Assaf, M., Naumann, D.A., Signoles, J., Totel, E., Tronel, F.: Hypercollecting semantics and its application to static analysis of information flow. In: *POPL* (2017)
6. Berzish, M., Ganesh, V., Zheng, Y.: Z3str3: A string solver with theory-aware heuristics. In: *FMCAD* (2017)
7. Blanchet, B., Cousot, P., Cousot, R., Feret, J., Mauborgne, L., Miné, A., Monniaux, D., Rival, X.: A static analyzer for large safety-critical software. In: *PLDI* (2003)
8. Bouajjani, A., Drăgoi, C., Enea, C., Rezine, A., Sighireanu, M.: Invariant synthesis for programs manipulating lists with unbounded data. In: *CAV* (2010)
9. Bouajjani, A., Drăgoi, C., Enea, C., Sighireanu, M.: Abstract domains for automated reasoning about list-manipulating programs with infinite data. In: *VMCAI* (2012)
10. Bouajjani, A., Drăgoi, C., Enea, C., Sighireanu, M.: Accurate invariant checking for programs manipulating lists and arrays with infinite data. In: *ATVA* (2012)
11. Bouajjani, A., Drăgoi, C., Enea, C., Sighireanu, M.: A logic-based framework for reasoning about composite data structures. In: *CONCUR* (2009)
12. Bouajjani, A., Drăgoi, C., Enea, C., Sighireanu, M.: On inter-procedural analysis of programs with lists and data. In: *PLDI* (2011)
13. Calcagno, C., Distefano, D., O’Hearn, P., Yang, H.: Footprint analysis : A shape analysis that discovers preconditions. In: *SAS* (2007)
14. Calcagno, C., Distefano, D., O’Hearn, P., Yang, H.: Compositional shape analysis by means of bi-abduction. In: *POPL* (2009)
15. Chang, B.Y.E., Rival, X.: Relational inductive shape analysis. In: *POPL*. *ACM* (2008)
16. Chang, B.E., Dragoi, C., Manevich, R., Rinetzky, N., Rival, X.: Shape analysis. *FNT* (1-2) (2020)
17. Charguéraud, A.: Characteristic formulae for the verification of imperative programs. In: *ICFP* (2011)
18. Cousot, P., Cousot, R.: Abstract interpretation: A unified lattice model for static analysis of programs by construction or approximation of fixpoints. In: *POPL*. *ACM* (1977)
19. Cousot, P., Cousot, R.: Systematic design of program analysis frameworks. In: *POPL* (1979)
20. Cousot, P., Cousot, R., Logozzo, F.: A parametric segmentation functor for fully automatic and scalable array content analysis. In: *POPL* (2011)
21. Cousot, P., Halbwachs, N.: Automatic discovery of linear restraints among variables of a program. In: *POPL* (1978)

22. Cox, A., Chang, B.Y.E., Rival, X.: Automatic analysis of open objects in dynamic language programs. In: SAS (2014)
23. Curry, C., Le, Q.L.: Bi-abduction for shapes with ordered data (2020), arXiv, <https://arxiv.org/abs/2006.10439>
24. Darnis, N.: The generic data-structure library (2004), <https://directory.fsf.org/wiki/GDSL>
25. Distefano, D., Fähndrich, M., Logozzo, F., O’Hearn, P.: Scaling static analyses at facebook. CACM (2019)
26. Dor, N., Rodeh, M., Sagiv, S.: Csvg: towards a realistic tool for statically detecting all buffer overflows in c. In: PLDI (2003)
27. Ferrara, P., Burato, E., Spoto, F.: Security analysis of the OWASP benchmark with julia. In: ITASEC (2017)
28. Fiedor, T., Holík, L., Rogalewicz, A., Sinn, M., Vojnar, T., Zuleger, F.: From shapes to amortized complexity. In: VMCAI (2018)
29. Gopan, D., Reps, T.W., Sagiv, S.: A framework for numeric analysis of array operations. In: POPL (2005)
30. Gulwani, S., McCloskey, B., Tiwari, A.: Lifting abstract interpreters to quantified logical domains. In: POPL (2008)
31. Halbwachs, N., Péron, M.: Discovering properties about arrays in simple programs. In: PLDI (2008)
32. Holík, L., Lengál, O., Rogalewicz, A., Simáček, J., Vojnar, T.: Fully automated shape analysis based on forest automata. In: CAV (2013)
33. Holík, L., Peringer, P., Rogalewicz, A., Šoková, V., Vojnar, T., Zuleger, F.: Low-level bi-abduction. In: ECOOP (2022)
34. Inc., A.: The freertos kernel (2022), <https://github.com/FreeRTOS>
35. Iosif, R., Rogalewicz, A., Vojnar, T.: Deciding entailments in inductive separation logic with tree automata. In: ATVA (2014)
36. Jeannet, B., Miné, A.: Apron: A library of numerical abstract domains for static analysis. In: CAV (2009)
37. Jung, R., Krebbers, R., Jourdan, J.H., Bizjack, A., Birkedal, L., Dreyer, D.: Iris from the ground up: A modular foundation for higher-order concurrent separation logic. *Journal of Functional Programming* (2018)
38. Kiezun, A., Ganesh, V., Artzi, S., Guo, P.J., Hooimeijer, P., Ernst, M.D.: Hampi: A solver for word equations over strings, regular expressions, and context-free grammars. *ACM Transaction of Software Engineering Methodology* (2013)
39. Le, Q.L., Gherghina, C., Qin, S., Chin, W.N.: Shape analysis via second-order bi-abduction. In: CAV (2014)
40. Li, H., Berenger, F., Chang, B.Y.E., Rival, X.: Semantic-directed clumping of disjunctive abstract states. In: POPL (2017)
41. Li, H., Rival, X., Chang, B.E.: Shape analysis for unstructured sharing. In: SAS (2015)
42. Liang, T., Reynolds, A., Tinelli, C., Barrett, C., Deters, M.: A dpll(t) theory solver for a theory of strings and regular expressions. In: CAV (2014)
43. Liu, J., Chen, L., Rival, X.: Automatic verification of embedded system code manipulating dynamic structures stored in contiguous regions. *IEEE Transactions on Computer Aided Design and Integration of Circuits Systems* (2018)
44. Loginov, A., Reps, T., Sagiv, M.: Automated verification of the deutsch-schorr-waite tree-traversal algorithm. In: SAS (2006)
45. Magill, S., Tsai, M.H., Lee, P., Tsay, Y.K.: Automatic numeric abstractions for heap-manipulating programs. In: POPL (2010)

46. Makanin, G.S.: The problem of solvability of equations in a free semigroup. *Mathematics of the USSR—Sbornik* **32**(4) (1977)
47. Midtgaard, J., Nielson, F., Nielson, H.R.: A parametric abstract domain for lattice-valued regular expressions. In: SAS (2016)
48. Miné, A.: The octagon abstract domain. HOSC (2006)
49. Negrini, L., Arceri, V., Ferrara, P., Cortesi, A.: Twinning automata and regular expressions for string static analysis. In: VMCAI (2021)
50. O’Hearn, P.: Resources, concurrency and local reasoning. In: CONCUR (2004)
51. Quine, W.V.: Concatenation as a basis for arithmetic. *Journal of Symbolic Logic* **11**(4) (1946). <https://doi.org/10.2307/2268308>
52. Reynolds, J.: Separation logic: A logic for shared mutable data structures. In: LICS (2002)
53. Sagiv, M., Reps, T., Wilhelm, R.: Solving shape-analysis problems in languages with destructive updating. TOPLAS (1998)
54. Singh, G., Püschel, M., Vechev, M.T.: Fast polyhedra abstract domain. In: POPL (2017)
55. Ta, Q.T., Le, T.C., Khoo, S.C., Chin, W.N.: Automated mutual explicit induction proof in separation logic. In: FAC (2016)
56. Torvalds, L.: The linux kernel (2022), <https://git.kernel.org>
57. Urban, C.: The abstract domain of segmented ranking functions. In: SAS (2013)
58. Vafeiadis, V.: Shape-value abstraction for verifying linearizability. In: VMCAI (2009)
59. Vafeiadis, V., Parkinson, M.: A marriage of rely/guarantee and separation logic. In: CONCUR (2007)



## A Experimental results

Example	without seq			with seq parameters					property verified
	time	#iter	property verified	time				#iter	
	all			all	num	seq	shape		
Linux lists									
Init	1.1		✓	2.6	0.2	0.3	1.1		✓
Input	13.6		✓	21.4	2.7	2.4	8.2		✓
Output	22.7		✓	31.5	4.8	4.8	10.5		✓
Output (empty)	33.8		✓	9.3	1.4	1.0	2.5		✓
FreeRTOS lists									
vListInit	4.3		✓	6.1	1.3	0.4	0.6		✓
vListInsertEnd	23.8		✓	40.3	10.8	1.8	5.3		✓
vListInsert	87.4	4	✓	370.5	202.4	27.2	37.9	4	✓
vListRemove	47.5		✓	163.4	82.6	9.2	20.0		✓
GDSDL (lists)									
Alloc	12.0		✓	14.8	2.2	1.4	3.0		✓
Flush	24.3	2	✓	59.4	18.4	5.4	16.1	2	✓
Free	35.3	2	✓(†)	79.9	25.1	7.4	24.0	2	✓(†)
Get size	3.6		✓	6.1	2.5	0.3	0.9		✓
Is empty (non-empty)	8.1		✓(†)	14.0	4.6	1.2	3.7		✓(†)
Is empty (empty)	8.1		✓(†)	24.3	12.1	1.8	3.9		✓(†)
Get head (empty)	9.5		✓	14.7	4.5	1.1	3.6		✓
Get head (non-empty)	9.4		✓	29.8	14.5	1.8	4.3		✓
Get tail (empty)	11.1		✓(†)	26.6	14.4	1.6	4.7		✓(†)
Get tail (non-empty)	11.0		✓(†)	58.5	35.5	2.8	6.6		✓(†)
Insert head	23.2		✓	42.9	13.3	2.9	10.5		✓
Insert tail	25.0		✓(†)	54.3	20.5	3.5	11.8		✓(†)
Remove head (empty)	34.1		✓	111.9	50.9	6.5	25.4		✓
Remove head (non-empty)	34.0		✓	16.3	5.7	1.1	3.7		✓
Remove tail (empty)	49.5		✓	284.8	165.0	13.6	39.3		✓
Remove tail (non-empty)	49.5		✓	16.2	5.7	1.1	3.6		✓
Search max	69.7	5	✓	708.4	429.7	43.1	145.7	5	PrSafe Fc
Search min	69.4	5	✓	634.0	380.3	35.4	131.2	5	PrSafe Fc
Search by position	104.5	3;2	✗(†)	1182.8	796.3	40.7	108.2	3;3	✓(†)

**Table 3.** Experimental results (Time in milliseconds averaged over 100 runs. For loop iterations, disjoint loops are separated by a semicolon, nested loops by a comma, and the first number corresponds to the outer loop. For inner loops, we take the maximum number of iterations needed to stabilize it.)

## B Detailed presentation of sequence operators

$$\frac{\text{len}_S = 0}{S = []} \quad \frac{\text{min}_S > \text{max}_S}{S = []} \quad \frac{S = [] \quad S = E \quad S' \in \mathbf{fv}(E) \cap \mathbb{X}_s}{S' = []}$$

$$\frac{S = E \quad \mathbf{fv}(E) \cap \mathbb{X}_n = \emptyset \quad \forall S' \in \mathbf{fv}(E), S' = []}{S = []}$$

(a) Inference rules for empty sequence

$$\frac{[\alpha].E_1 = [\beta].E_2 \quad \alpha = \beta \quad E_1 = E_2}{S.E = E' \quad S = []} \quad \frac{S_1.E_1 = S_2.E_2 \quad \text{len}_{S_1} = \text{len}_{S_2} \quad S_1 = S_2 \quad E_1 = E_2}{E = E'}$$

(b) Inference rules to decompose equality constraints

$$\frac{S = E \quad S' \in \mathbf{fv}(E) \cap \mathbb{X}_s}{\text{min}_S \leq \text{min}_{S'} \quad \text{max}_{S'} \leq \text{max}_S} \quad \frac{S = [] \quad \alpha \in \mathbb{X}_n}{\text{max}_S \leq \alpha \leq \text{min}_S}$$

$$\frac{S = [\alpha]}{\text{min}_S = \alpha = \text{max}_S} \quad \frac{S = E \quad \delta \in \mathbf{fv}(E) \cap \mathbb{X}_n}{\text{min}_S \leq \delta \leq \text{max}_S}$$

(c) Inference rules for bound constraints (general case)

$$\frac{S = \mathbf{sort}(S) \quad S = \dots S' \dots S'' \dots}{\text{max}_{S'} \leq \text{min}_{S''}}$$

$$\frac{S = \mathbf{sort}(S) \quad S = \dots [\alpha] \dots S' \dots}{\alpha \leq \text{min}_{S'}} \quad \frac{S = \mathbf{sort}(S) \quad S = \dots [\alpha] \dots [\beta] \dots}{\alpha \leq \beta}$$

$$\frac{S = \mathbf{sort}(S) \quad S = [\alpha] \dots}{\alpha = \text{min}_S}$$

(d) Inference rules for bound constraints (sorted case)

$$\frac{S = []}{S = \mathbf{sort}(S)} \quad \frac{S = \mathbf{sort}(E)}{S = \mathbf{sort}(S)} \quad \frac{S = \mathbf{sort}(S) \quad S = \dots S' \dots}{S' = \mathbf{sort}(S')}$$

$$S = S_{1,1} \dots S_{1,l_1} \cdot [\delta_1] \cdot S_{2,1} \dots S_{k-1,l_{k-1}} \cdot [\delta_{k-1}] \cdot S_{k,1} \dots S_{k,l_k}$$

$$\forall i, \delta_i \leq \delta_{i+1} \quad \forall i \in [1, k-1], \forall j, j', \text{max}_{S_{i,j}} \leq \delta_i \leq \text{min}_{S_{i+1,j'}}$$

$$\forall i, \forall j, S_{i,j} = \mathbf{sort}(S_{i,j}) \quad \forall i, \forall j < j', \text{max}_{S_{i,j}} \leq \text{min}_{S_{i,j'}}$$

$$\frac{}{S = \mathbf{sort}(S)}$$

(e) Inference rules for sortedness

**Fig. 12.** Inference rules for guard in sequence domain

## C Proof of concatenation lemma

*Proof.* (Concatenation lemma for list, Lemma 1) Let  $(m_0, \sigma)$  be a concrete state in the concretization of  $(m_0^\#, \sigma_m^\#)$ . This means that there exists  $m', m''$ , two disjoint memory such that:  $m_1, \sigma \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_0, \alpha_1, S_1 \square)$ ,  $m_2, \sigma \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_1, \alpha_2, S_2 \square)$ , and  $m_1 \uplus m_2 = m_0$ .

We prove that  $m_1 \uplus m_1, \sigma \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_0, \alpha_2, S \square)$  by structural induction on the judgments used to establish that  $m_1, \sigma \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_0, \alpha_1, S_1 \square)$ .

*Empty case.* In that situation,  $\sigma_n(\alpha_0) = \sigma_n(\alpha_1)$ ,  $\sigma_s(S_1) = \sigma_s(S)$ , and  $m_0 = m_2$ . By hypothesis on  $m_2$ , we establish that  $m_0, \sigma \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_0, \alpha_2, S \square)$

*Non-empty case.* From the premise of this rules it follows that there exist two disjoint concrete memories  $m'_1, m''_1$ , and symbolic variables  $\alpha'_0 \delta, S'_1$  such that the following hold:

$$\begin{aligned} m_1 &= m'_1 \uplus m''_1 \\ m'_1, \sigma &\models_{\mathbb{M}} \mathbf{lseg}_s(\alpha'_0, \alpha_1, S'_1 \square) \\ m''_1, \sigma &\models_{\mathbb{M}} \alpha_1 \mathbf{next} \mapsto \alpha'_1 * \alpha_1 \mathbf{.data} \mapsto \delta \\ \sigma &\models_s S_1 = [\delta].S'_1 \end{aligned}$$

Let's define  $S' := S'_1.S_2$ , it follows that  $\sigma \models_s S = [\delta].S'$ . By induction hypothesis, we establish that  $m'_1 \uplus m_2, \sigma \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha'_1, \alpha_3, S' \square)$ . Finally, from the premise on  $m''_1$  and the latter, we apply the non-empty list segment rule to conclude that:  $m\sigma \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_1, \alpha_3, S \square)$

Furthermore, let  $(m, \sigma_1)$  be a concrete state in  $\gamma_{\mathbb{S}}(m_1^\#, \sigma^\#)$ . We prove by induction on  $|\sigma_1(S_1)|$  that there exists  $\sigma_0$  such that  $\forall \beta \neq \alpha_1, \sigma_0(\beta) = \sigma_1(\beta)$ , and  $(m, \sigma_0) \in \gamma_{\mathbb{S}}(m_0^\#, \sigma^\#)$ .

*Null case* If  $\sigma_1(S_1) = \varepsilon$ , this means that  $\sigma_1 \models_s S = S_2$ . So we simply define  $\sigma_0 := \sigma_1\{\alpha_1 \mapsto \sigma_1(\alpha_0)\}$ , and we have:

- $\sigma_0 \models_s S = S_1.S_2$ ,
- $\emptyset, \sigma_0 \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_0, \alpha_1, S_1 \square)$ ,
- $m, \sigma_0 \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_1, \alpha_2, S_2 \square)$ ,

Therefore, we conclude that:  $(\sigma_0, mem) \in \gamma_{\mathbb{S}}(m_0^\#, \sigma^\#)$ .

*Successor case* If  $|\sigma_1(S_1)| > 0$ , then  $\sigma_1(S) \neq \varepsilon$ , and we know that the only possible rule used to establish that  $m, \sigma_1 \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_0, \alpha_2, S \square)$  is the non-empty one. So we deduce that there exist  $m', m''$  such that for some fresh variables  $S', \alpha'_0, \delta$ :

- $m', \sigma_1 \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha'_0, \alpha_2, S' \square)$
- $m'', \sigma_1 \models_{\mathbb{M}} \alpha_0 \mathbf{.next} \mapsto \alpha'_0 * \alpha_0 \mathbf{.data} \mapsto \delta$
- $\sigma_1 \models_s [\delta].S'$
- $m = m' \uplus m''$ .

We define  $\sigma'_1 := \sigma_1 S'_1 \mapsto \sigma_1(S_1)_{[1, |\sigma_1(S_1)|]}$ . That is to say  $\sigma'_1$  maps to some fresh auxiliary sequence variable  $S'_1$  the valuation of  $S_1$ , without the first character. Note that  $\sigma'_1$  differs from  $\sigma_1$  only when evaluated at  $S'_1$  (and  $\mathbf{min}_{S_1}, \dots$ ). The following holds:

$$\begin{aligned} (m', \sigma'_1) &\models_{\mathbb{M}} \mathbf{lseg}_s(\alpha'_0, \alpha_2, S' \square) \\ \sigma'_1 &\models_s S' = S'_1.S_2 \end{aligned}$$

So, we can apply the induction hypothesis: We know there exists some  $\sigma'_0$  such that it only differs from  $\sigma'_1$  when evaluated on  $\alpha_1$  and such that  $m', \sigma'_0 \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha'_0, \alpha_1, S'_1 \square) * \mathbf{lseg}_s(\alpha_1, \alpha_2, S_2 \square)$ . This means that there is two disjoint concrete memories  $m'_1, m'_2$  such that

- $m' = m'_1 \uplus m'_2$
- $m'_1, \sigma'_0 \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha'_0, \alpha_1, S'_1 \square)$
- $m'_2, \sigma'_0 \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_1, \alpha_2, S_2 \square)$

Therefore, we can establish using the non-empty segment rule that:  $m'' \uplus m'_1, \sigma'_0 \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_0, \alpha_1, S_1 \square)$ . If, we define  $\sigma_0 := \sigma'_0 \{S'_1 \mapsto cs_1(S'_1)\}$ , the latter still holds since  $S'_1$  does not occur free in the abstract memory state, and we also know that  $m'_2, \sigma_0 \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_1, \alpha_2, S_2 \square)$ . And  $\sigma_0$  only differs from  $\sigma_1$  when evaluated on  $\alpha_1$ . And since  $m = m'' \uplus m' = m'' \uplus m'_1 \uplus m'_2$ , we infer that  $m, \sigma_0 \models_{\mathbb{M}} m_0^\sharp$ . To conclude we have:  $m, \sigma_0 \in \gamma_{\mathbb{S}}(m_0^\sharp, \sigma^\sharp)$ .

The proof of Lemma 2 is established in a similar manner, by induction on the first segment.

$$\begin{array}{c}
 \frac{\sigma_n(\alpha_0) = \sigma_n(\alpha_1) \quad \sigma_n, \sigma_s \models_s S = \square}{\emptyset, (\sigma_n, \sigma_m, \sigma_s) \models_{\mathbb{M}} \mathbf{tree}_s(\alpha, S)} \\
 \\
 \frac{m, (\sigma_n, \sigma_m, \sigma_s) \models_{\mathbb{M}} \begin{array}{l} \alpha.l \mapsto \alpha_l * \alpha.r \mapsto \alpha_r * \alpha.d \mapsto \alpha_d \\ * \mathbf{tree}_s(\alpha_l, S_l) * \mathbf{tree}_s(\alpha_r, S_r) \end{array} \quad \sigma_n(\alpha) \neq 0 \quad \sigma_n, \sigma_s \models_s S = S_l.[\alpha_d].S_r \quad S_l, S_r \text{ fresh}}{m, (\sigma_n, \sigma_m, \sigma_s) \models_{\mathbb{M}} \mathbf{tree}_s(\alpha, S)} \\
 \\
 \frac{\sigma_n(\alpha_0) = \sigma_n(\alpha_1) \quad \sigma_n, \sigma_s \models_s S_l = S_r = \square}{\emptyset, (\sigma_n, \sigma_m, \sigma_s) \models_{\mathbb{M}} \mathbf{tree}_s(\alpha_0, \alpha_1, S_l \square S_r)} \\
 \\
 \frac{m, (\sigma_n, \sigma_m, \sigma_s) \models_{\mathbb{M}} \begin{array}{l} \alpha.l \mapsto \alpha_l * \alpha.r \mapsto \alpha_r * \alpha.d \mapsto \alpha_d \\ * \mathbf{tree}_s(\alpha_l, S'_l) * \mathbf{tree}_s(\alpha_r, \alpha_1, S'_l \square S'_r) \end{array} \quad \sigma_n(\alpha) \neq 0 \quad \sigma_n, \sigma_s \models_s S_l = S'_l.[\alpha_d].S'_l \quad \sigma_n, \sigma_s \models_s S_r = S'_r \quad S', S'_l, S'_r \text{ fresh}}{m, (\sigma_n, \sigma_m, \sigma_s) \models_{\mathbb{M}} \mathbf{tree}_s(\alpha_0, \alpha_1, S_l \square S_r)} \\
 \\
 \frac{m, (\sigma_n, \sigma_m, \sigma_s) \models_{\mathbb{M}} \begin{array}{l} \alpha.l \mapsto \alpha_l * \alpha.r \mapsto \alpha_r * \alpha.d \mapsto \alpha_d \\ * \mathbf{tree}_s(\alpha_l, \alpha_1, S'_l \square S'_r) * \mathbf{tree}_s(\alpha_r, S') \end{array} \quad \sigma_n(\alpha) \neq 0 \quad \sigma_n, \sigma_s \models_s S_l = S'_l \quad \sigma_n, \sigma_s \models_s S_r = S'_r.[\alpha_d].S' \quad S', S'_l, S'_r \text{ fresh}}{m, (\sigma_n, \sigma_m, \sigma_s) \models_{\mathbb{M}} \mathbf{tree}_s(\alpha_0, \alpha_1, S_l \square S_r)}
 \end{array}$$

**Fig. 13.** Concretization of tree memory predicates

*Proof.* (Concatenation lemma for binary trees, Lemma 2) We here prove the segment/full predicate concatenation. Let  $(m_0, \sigma)$  be an abstract state in the concretization of  $(m_0^\sharp, \sigma_m^\sharp)$ . Therefore, we know that there exists two disjoint concrete memories  $m_1$  and  $m_2$  such that  $m_1, \sigma \models_{\mathbb{M}} \mathbf{tree}_s(\alpha, \alpha', S'_l \square S'_r)$ ,  $m_2, \sigma \models_{\mathbb{M}} \mathbf{tree}_s(\alpha', S')$ , and  $m_1 \uplus m_2 = m_0$ .

We prove that  $m_1 \uplus m_1, \sigma \models_{\mathbb{M}} \mathbf{tree}_s(\alpha_0, S)$  by structural induction on the judgments used to establish that  $m_1, \sigma \models_{\mathbb{M}} \mathbf{lseg}_s(\alpha_0, \alpha_1, S_1 \square)$ .

*Empty case* In that situation,  $\sigma_n(\alpha) = \sigma_n(\alpha')$ ,  $\sigma_s(S') = \sigma_s(S)$ , and  $m_0 = m_2$ . By hypothesis on  $m_2$ , we establish that  $m_0, \sigma \models_{\mathbb{M}} \mathbf{tree}_s(\alpha, S)$

*Left case* From the premise of this rules it follows that there exist two disjoint concrete memories  $m'_1, m''_1$ , and symbolic variables  $\alpha_l, \alpha_r, \alpha_d, S'_l, S''_l; S'_r, S''_r$  such that the following hold:

$$\begin{aligned} m_1 &= m'_1 \uplus m''_1 \\ m'_1, \sigma &\models_{\mathbb{M}} \mathbf{treeseq}_s(\alpha_l, \alpha', S'_l \square S''_l) \\ m''_1, \sigma &\models_{\mathbb{M}} \alpha.l \mapsto \alpha_l * \alpha.r \mapsto \alpha_r * \alpha.d \mapsto \alpha_d * \mathbf{tree}_s(\alpha_r, S'') \\ \sigma &\models_s S'_l = S''_l \\ \sigma &\models_s S'_r = S''_r.[\alpha_d].S'' \end{aligned}$$

Let's define  $S_0 := S'_l.S'.S''_r$ , it follows that  $\sigma \models_s S = S_0.[\alpha_d].S''$ . By induction hypothesis, we establish that  $m'_1 \uplus m_2, \sigma \models_{\mathbb{M}} \mathbf{tree}_s(\alpha_l, \alpha', S'')$ . Finally, from the premise on  $m''_1$  and the latter, we apply the non-empty list segment rule to conclude that:  $m\sigma \models_{\mathbb{M}} \mathbf{tree}_s(\alpha, S)$

The *right case* is proved similarly.

## D Proofs of sequence abstract operators

*Proof.* (Soundness of  $\mathbf{guard}_{\Sigma}$ , Theorem 1)

**Lemma 3 (Soundness of definition inlining).** *For any definition constraints  $S = E$  and  $S' = E'$ , and any concrete state  $(\sigma_n, \sigma_s)$ , if  $\sigma_n, \sigma_s \models_s S' = E'$ , then  $\sigma_n, \sigma_s \models_s S = E$  if and only if  $\sigma_n, \sigma_s \models S = E[S' \leftarrow E']$ .*

*Proof.* This is proved by a simple induction on  $E$ .

**Lemma 4 (Length constraints saturation).** *For any concrete state  $(\sigma_n, \sigma_m, \sigma_s) \in \Sigma$ , any definition constraint  $S = E$ , if  $\sigma_n, \sigma_s \models_s S = E$ , then  $\sigma_n \models_n \mathbf{len}_S = \tau_{\mathbf{len}}(E)$ .*

*Proof.* First, we prove by induction on  $E$  that if  $\llbracket E \rrbracket_s(\sigma_n, \sigma_s) = a_1 \dots a_k$ , then  $\llbracket \tau_{\mathbf{len}}(E) \rrbracket_n(\sigma_n) = k$ . The only interesting case is when  $E$  is only a sequence variable  $S'$ . We simply use the hypothesis of  $(\sigma_n, \sigma_m, \sigma_s) \in \Sigma$  to conclude that  $\sigma_n(\mathbf{len}_{S'}) = |\sigma_s(S')|$ .

Finally, since  $\sigma_n, \sigma_s \models_s S = E$ , we have that  $\sigma_s(S) = \llbracket E \rrbracket_s(\sigma_n, \sigma_s) = a_1 \dots a_k$ , and because the concrete state is well-formed, we conclude that  $\sigma_n(\mathbf{len}_S) = k = \llbracket \tau_{\mathbf{len}} \rrbracket_n(\sigma_n)$ .

The soundness of multi-set constraints saturation is proved in a similar manner.

**Lemma 5 (Soundness of empty sequence variables detection).** *The rules presented in figure 12(a) are sound. That is to say, for any concrete state  $(\sigma_n, \sigma_m, \sigma_s)$  and a rule whose premises (resp. conclusions) are  $P_i, \dots, P_k$  (resp.  $C_1, \dots, C_l$ ), if  $\forall i, (\sigma_n, \sigma_m, \sigma_s) \models P_i$ , then  $\forall j, (\sigma_n, \sigma_m, \sigma_s) \models C_j$ .*

*Proof.* The first rule is a direct consequence of the fact that for any concrete state  $(\sigma_n, \sigma_m, \sigma_s) \in \Sigma$ ,  $\sigma_s(S) = \varepsilon \Leftrightarrow \sigma_n(\mathbf{len}_S) = 0$ .

To establish the second one, we observe that  $\sigma_s(S) = a_1 \dots a_n \neq \varepsilon \Rightarrow \sigma_n(\mathbf{min}_S) \leq a_1 \leq \sigma_n(\mathbf{max}_S)$ . The contraposition gives the expected result.

For the third rule, let's assume there is a sequence variable  $S' \in \mathbf{fv}(E) \cap \mathbb{X}_s$  such that  $\sigma_s(S) \neq \varepsilon$ . By induction on  $E$  we prove that  $\llbracket E \rrbracket_s(\sigma_n, \sigma_s) \neq \varepsilon$ . This contradicts the hypothesis  $S = []$ .

The fourth rule is demonstrated by induction on  $E$ . The sequence variable case is proved using the assumption that all free-variables of  $E$  are empty. The value symbolic variable case is impossible since  $E$  does not have such free variables. The other three cases (empty sequence, **sort** and concatenation) are straightforward.

**Lemma 6 (Soundness of decomposition of equality constraints).** *All rules presented in figure 12(b) are sound.*

*Proof.* The first two rules, are a result of the fact same-length prefixes of a word are equals. For the third rule, this is a consequence of lemma 3.

**Lemma 7 (Soundness of bound saturation rules).** *All rules described in figures 12(c) and (d) are sound.*

*Proof.* For the first rule, the case where  $\sigma_s(S') = \varepsilon$  is straightforward. Now consider the case that  $\sigma_s(S')$  is non-empty, then so is  $\sigma_s(S)$ . And since any element of  $\sigma_s(S')$  is also an element of  $\sigma_s(S)$ , the inequalities on extremal values hold.

The next three rules are straightforward.

*Remark 2.* We emphasize that stating that for any sequence variable  $S$ ,  $\mathbf{min}_S \leq \mathbf{max}_S$  is unsound, since we have by convention that  $S = [] \Rightarrow \mathbf{min}_S = +\infty$  and  $\mathbf{max}_S = -\infty$ . Generally bounding from below  $\mathbf{max}_S$  and from above  $\mathbf{min}_S$  should be done with extreme care. That's why the two rules doing so are the ones implying that  $S$  is non-empty.

**Lemma 8 (Soundness of inference rule for sortedness).** *All rules presented in figure 12(e) are sound.*

*Proof.* The first two rules are straightforward.

The third one is established by induction on  $E$  where  $E$  is the definition of  $S$  where  $S'$  occurs. Notice that the premise  $S = \dots S' \dots$  forbids  $S'$  to occur inside a call to **sort**.

For the last rule, let's define  $\sigma_s(S) = a_1 \dots a_m$  and take  $n \in [1, m - 1]$ . We reason by case disjunction by considering all possible sources of  $a_n$  and  $a_{i+1}$ .

- If both  $a_n$  and  $a_{n+1}$  come from evaluation of value symbolic variables  $[\delta_i]$ ,  $[\delta_{i+1}]$ , then the second premise, ensures that  $a_n \preceq a_{n+1}$ . This case occurs when  $l_k = 0$ .
- If  $a_n$  comes from a value symbolic variable  $[\delta_i]$ , and  $a_{n+1}$ , from a sequence variable  $S_{i+1,j}$ , then the third hypothesis gives the expected result.
- The case where  $a_n$  comes from a sequence variable and  $a_{n+1}$ , from value symbolic variable, is proved likewise.

- Finally, let's consider the case where both  $a_n$  and  $a_{n+1}$  come from a sequence variable. It is impossible they come from sequence variables with distinct first indexes  $S_{i,j}, S_{i',j'}$  where  $i \neq i'$ , because there is necessarily a value between the evaluation of these sequence variables. Therefore,  $a_n$  and  $a_{n+1}$  come from  $S_{i,j}$  and  $S_{i,j'}$ , respectively.
  - If  $j = j'$ , then we conclude directly using the fourth premise.
  - Otherwise, thanks to the last premise, we infer:  $a_n \preceq \sigma_n(\max_{S_{i,j}}) \preceq \sigma_n(\min_{S_{i,j'}}) \preceq a_{n+1}$ .

*Remark 3.* As observed in Section 3 the final premise cannot be weakened by only considering the case where  $j' = j + 1$ . To provide a clear counter-example, consider the abstract value defined as  $S = S_1.S'.S_2 \wedge \max_{S_1} \leq \min_{S'} \wedge \max_{S'} \leq \max_{S_2}$ . These constraints are not sufficient to ensure that  $S$  is sorted. Indeed, one possible concrete state of this abstract value, where  $\sigma_s(S)$  is not sorted is:

$$\left\{ \begin{array}{ll} S \mapsto 31 & S_1 \mapsto 3 \\ S' \mapsto \varepsilon & S_2 \mapsto 1 \end{array} \right\}$$

*Proof (Soundness of  $\text{verif}_{\Sigma}$ , Theorem 2).* For the cases where  $\sigma_s^\#$  is  $\perp$  or the constraint is a bound constraint, the inclusion is a consequence of the definition of  $\gamma_{\Sigma}$ .

For definition constraints, the soundness is a result of lemma 3.

Let's prove the soundness of the following rule:

$$\frac{S = \mathbf{sort}(S) \quad \mathbf{multi}_S = \tau_{\mathbf{mul}}(E)}{S = \mathbf{sort}(E)}$$

Let  $(\sigma_n, \sigma_m, \sigma_s)$  be a well-formed concrete state, Since  $\mathbf{multi}_S = \tau_{\mathbf{mul}}(E)$  holds, there is a bijection  $\pi$  such that  $\sigma_s(S) = a_1 \dots a_k$  and  $\llbracket E \rrbracket_s = a_{\pi(1)} \dots a_{\pi(k)}$ . Moreover,  $S = \mathbf{sort}(S)$ , ensures that  $\forall i \in [1, k - 1], a_i \preceq a_{i+1}$ .