

Automatic Inference of Heap Properties Exploiting Value Domains

Pietro Ferrara¹, Peter Müller², and Milos Novacek²

¹ IBM Thomas J. Watson Research Center, USA

`pietroferrara@us.ibm.com`

² Department of Computer Science, ETH Zurich, Switzerland

`peter.mueller@inf.ethz.ch, milos.novacek@inf.ethz.ch`

Abstract. Effective static analyses of heap-manipulating programs need to track precise information about the heap structures and the values computed by the program. Most existing heap analyses rely on manual annotations to precisely analyze general and, in particular, recursive, heap structures. Moreover, they either do not exploit value information to obtain more precise heap information or require more annotations for this purpose. In this paper, we present a combined heap and value analysis that infers complex invariants for recursive heap structures such as lists and trees, including relations between value fields of heap-allocated objects. Our analysis uses a novel notion of edge-local identifiers to track value information about the source and target of a pointer, even if these are summary nodes. With each potential pointer in the heap, our analysis associates value information that describes in which states the pointer may exist, and uses this information to improve the precision of the analysis by pruning infeasible heap structures. Our analysis has been implemented in the static analyzer Sample; experimental results show that it can automatically infer invariants for data structures, for which state-of-the-art analyses require manual annotations.

1 Introduction

Effective static analyses of heap-manipulating programs need to track precise information about the heap structures and the values computed by a program. Heap and value information is not independent: heap information determines which locations need to be tracked by a value analysis, and information about value fields may be useful to obtain more precise heap information, for instance, to rule out certain forms of aliasing. Moreover, many interesting invariants of heap-manipulating programs combine heap and value information such as the invariant that a heap structure is a sorted linked list.

Despite these connections, heap and value analyses have often been treated as orthogonal problems. Some existing heap analyses such as TVLA [17] rely on manual instrumentation to infer invariants that combine heap and value information. However, TVLA does not support general value domains, which limits, for instance, arithmetical reasoning. Recent work addresses this issue by combining TVLA with value domains, but still requires the user to provide predicates

to track and exchange information between the heap and value domains [21], or is not able to track complex invariants over recursive data structures [13]. Chang and Rival [4] present an efficient inference for combined heap and value invariants, which also relies on user-provided predicates. Other analyses do not require manual annotations [1,2], but are specific to programs that manipulate certain data structures such as singly-linked lists.

In this paper, we present a combined heap and value analysis—expressed as an abstract interpretation [7]—that infers complex invariants of heap structures. It is automatic in the sense that it uses only the information included in the program, without relying on manual annotations. Our analysis uses a graph-based abstraction of heaps, where each edge in the graph represents a *potential* pointer in the concrete heap. Each edge is associated with an abstract value state that characterizes in which concrete states this pointer might actually exist. The value states on the edges allow our analysis to represent disjunctive information in a single heap graph (like the bracketing constraints in Dillig et al.’s Fluid Updates [9]). They are also used to improve the precision of the analysis when value information implies that certain pointer chains cannot exist in concrete heaps. Our analysis can be instantiated with different value domains to obtain different trade-offs between precision and efficiency.

Like many heap analyses, we use summary nodes to abstract over sets of concrete objects. A key innovation of our analysis is to introduce *edge-local identifiers* for the source and target of each edge in the heap graph. An edge-local identifier represents a field of one particular concrete object, even when the object is abstracted by a summary node. By having identifiers per edge, the value analysis may relate the fields of the source and the target of a concrete pointer and, thus, track inductive invariants such as the sortedness of a linked list.

Example. Method `increasingList` in Fig. 1 creates and returns a linked list. If parameter `v` is non-positive, the list is empty, that is, `result` is null (invariant I1). Otherwise, the list satisfies the following properties: it is non-empty, that is, the `result` is non-null (I2), the first node has value 0 (I3), the values of all other nodes are one larger than their predecessor’s (I4), and the value of the last node is `v − 1` (I5). Note that these invariants imply that the list is acyclic and has `v` nodes.

```

1 Node increasingList (int v) {
2   Node result = null;
3   int i = v;
4   while (i > 0) {
5     Node p = new Node();
6     p.next = result;
7     p.val = i - 1;
8     result = p;
9     i = i - 1;
10  }
11  return result;
12 }

```

Fig. 1: Running example.

Fig. 2 shows the abstract state that our analysis infers at the end of method `increasingList`. Here, we use a numerical domain such as Polyhedra [8] or Octagon [22] for the abstract states associated with each edge in the graph. The figure shows the relevant constraints from these states. They are expressed in terms of parameter `v` and the edge-local identifiers (Src, val) and (Trg, val) , which refer to the `val` field of the source and target of a pointer, respectively.

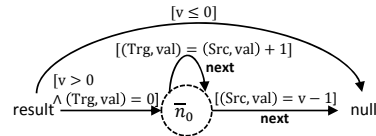


Fig. 2: The abstract heap state inferred at line 11 of Fig. 1.

The abstract state reflects the five invariants stated above. Variable **result** is null if the constraints on the corresponding edge hold, that is, if v is non-positive (I1). Otherwise, **result** points to the summary node \bar{n}_0 , which implies that it is non-null (I2). This example illustrates that our analysis represents disjunctive information in a single graph: both possible values of **result** are represented by the same graph, and we use value information to determine the states in which each pointer may exist. The constraint $(\text{Trg}, \text{val}) = 0$ on the edge from **result** to the summary node \bar{n}_0 expresses that the first list node has value 0 (I3). Note that the edge-local identifier allows us to express properties of a single object, even if it is abstracted by a summary node. The same feature is used in the constraint $(\text{Trg}, \text{val}) = (\text{Src}, \text{val}) + 1$ on the edge from \bar{n}_0 to itself to express invariant I4. Finally, the constraint $(\text{Src}, \text{val}) = v - 1$ on the edge from \bar{n}_0 to null expresses that the last list node has value $v - 1$ (I5). All five invariants are inferred automatically by our analysis without manual annotations.

Outline. Sec. 2 defines the language and the concrete domain. Sec. 3 formalizes the abstract domain, while Sec. 5 defines the abstract semantics. Sec. 6 reports the experimental results, Sec. 7 discusses related work, and Sec. 8 concludes.

2 Programming Language and Concrete Domain

We present our analysis for the small object-based language in Fig. 3. To simplify the formalization, we model local variables as fields of a special object \mathfrak{S} , that is, treat local variables as heap locations. We distinguish *reference field* and *value field* access expressions rAE and vAE , depending on the type of the accessed field. A reference expression **rexp** may be null, a reference field access expression, or an object creation. A value expression **vexp** may be a literal, a value field access expression, or a binary expression. Since the treatment of loops and conditionals is standard, the only relevant statements in **ST** are value and reference assignments.

$$\begin{aligned} \text{rAE} &::= \mathfrak{S}.f_r \mid \text{rAE}.f_r \\ \text{vAE} &::= \mathfrak{S}.f_v \mid \text{rAE}.f_v \\ \text{rexp} &::= \text{null} \mid \text{rAE} \mid \text{new } C \\ \text{vexp} &::= n \mid \text{vAE} \mid \text{vexp } \langle \text{op} \rangle \text{ vexp} \\ \text{op} &::= + \mid - \mid * \mid \dots \\ \text{ST} &::= \text{rAE} = \text{rexp} \mid \text{vAE} = \text{vexp} \end{aligned}$$

Fig. 3: Expressions and statements.

In the concrete domain, we partition the content of heap locations into values and references. Let Ref be the set of concrete references (objects and null), with $\mathfrak{S}, \text{null} \in \text{Ref}$, and let Val be the set of values. Let $\text{Field}_{\text{Ref}}$ and $\text{Field}_{\text{Val}}$ be finite sets of reference and value fields, respectively. An execution state consists of a *value store* and a *reference store*. We model a value store as a partial map in $\text{Store}_{\text{Val}} = (\text{Ref} \setminus \{\text{null}\}) \times \text{Field}_{\text{Val}} \rightharpoonup \text{Val}$ and a reference store as a partial map in $\text{Store}_{\text{Ref}} = (\text{Ref} \setminus \{\text{null}\}) \times \text{Field}_{\text{Ref}} \rightharpoonup \text{Ref}$. For each reference in their domain, these maps contain an entry for every field. We will refer to entries in a reference store as *concrete edges*. We define the set of all concrete states (*concrete heaps*) as $\Sigma = \text{Store}_{\text{Ref}} \times \text{Store}_{\text{Val}}$.

3 Abstract Domain and Operators

In this section, we present the abstract domain, the concretization function, as well as join and widening operators.

3.1 Abstract Domain

Let $\overline{\text{Ref}}$ be the set of *abstract references* (or *abstract nodes*) with $\mathfrak{I}, \text{null} \in \overline{\text{Ref}}$ (that is, we overload the symbols \mathfrak{I} and null to denote both concrete and abstract references). Each abstract node $\bar{n} \in \overline{\text{Ref}}$ represents either a single concrete non-null reference (*definite node*), or a non-empty set of concrete non-null references (*summary node*) with \mathfrak{I} and null being definite nodes. The functions in $\text{IsSummary} = \overline{\text{Ref}} \rightarrow \{\text{true}, \text{false}\}$ define whether a node is a summary node.

An *abstract reference store* in $\overline{\text{Store}}_{\overline{\text{Ref}}} = \mathcal{P}((\overline{\text{Ref}} \setminus \{\text{null}\}) \times \text{Field}_{\overline{\text{Ref}}} \times \overline{\text{Ref}})$ represents possible pointers between abstract nodes through reference fields. It can be interpreted as a directed graph where edges are labeled with a field name. Hence, we will refer to members of the abstract reference store as *abstract edges*. For an abstract edge $\bar{n}_1 \xrightarrow{f_r} \bar{n}_2$, we will refer to \bar{n}_1 as the *source* and to \bar{n}_2 as the *target* of the edge.

Our heap analysis is parameterized by an *abstract value domain* $\overline{\text{V}}$, which tracks information about value fields, for instance, relations among numerical values. Each abstract edge is associated with an abstract value state (*abstract condition*) via a map in $\overline{\text{Cond}} = (\overline{\text{Ref}} \setminus \{\text{null}\}) \times \text{Field}_{\overline{\text{Ref}}} \times \overline{\text{Ref}} \rightarrow \overline{\text{V}}$. The abstract condition of an abstract edge approximates the concrete value stores in which the edge exists. That is, our abstract domain tracks disjunctive information by having several edges with the same source and field, and associating them with different abstract conditions.

Abstract value states in $\overline{\text{V}}$ refer to memory locations via *abstract identifiers* $\overline{\text{ID}} = \overline{\text{Loc}} \cup \overline{\text{Eld}}$ where $\overline{\text{Loc}} = (\overline{\text{Ref}} \setminus \{\text{null}\}) \times \text{Field}_{\text{Val}}$ and $\overline{\text{Eld}} = \{\text{Src}, \text{Trg}\} \times \text{Field}_{\text{Val}}$. An identifier $(\bar{n}, f_v) \in \overline{\text{Loc}}$ represents the value field f_v of the concrete references abstracted by the node \bar{n} . *Edge-local identifiers* $(\text{Src}, f_v), (\text{Trg}, f_v) \in \overline{\text{Eld}}$ represent the value field f_v of the *single* concrete *source* or *target* reference of the concrete edges represented by an abstract edge. They track relations between the value fields of adjacent references in concrete heaps, which allows us to infer precise invariants on summary nodes. For instance, the constraint $(\text{Src}, \text{val}) \leq (\text{Trg}, \text{val})$ in the abstract condition of an abstract edge $\bar{n} \xrightarrow{\text{next}} \bar{n}$ expresses sortedness of the concrete list that is abstracted by the summary node \bar{n} .

We define the set of all abstract states (*abstract heaps*) as $\overline{\Sigma} = \overline{\text{Store}}_{\overline{\text{Ref}}} \times \overline{\text{Cond}} \times \text{IsSummary}$.

Example. The abstract heap in Fig. 4 depicts the loop invariant of the program in Fig. 1. Many of the constraints are similar to the constraints in Fig. 2. In particular, combining the abstract heap for the loop invariant with the negation of the loop guard (that is, $i \leq 0$) yields the information reflected in Fig. 2, for instance, that **result** is null iff $v \leq 0$ and that the first list node has value 0.

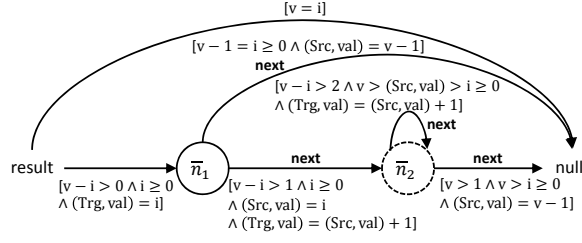


Fig. 4: The abstract heap representing the loop invariant at line 4 of the example in Fig. 1. Solid and dashed circles denote definite and summary nodes, respectively. Arrows depict abstract edges and are annotated with relevant constraints from their abstract conditions. To improve readability, we depict local reference variables as nodes and use local variables as identifiers in constraints, although the analysis models them as fields of \mathfrak{S} .

3.2 Concretization

In this section, we define the concretization function $\gamma : \overline{\Sigma} \rightarrow \mathcal{P}(\Sigma)$ that yields the set of concrete heaps represented by a given abstract heap.

We assume that our heap analysis is instantiated with a sound value analysis. Its concretization function $\gamma_{\overline{\mathbf{V}}} : \overline{\mathbf{V}} \rightarrow \mathcal{P}(\mathbf{ID} \rightarrow \mathcal{P}(\mathbf{Val}))$ yields a set of maps from abstract identifiers to sets of concrete values. These maps yield sets of concrete values rather than single values since an abstract value state may contain identifiers for fields of summary nodes, and the value analysis alone cannot concretize them. Let $references(st_{\text{Ref}})$ be the set of concrete references of a given concrete reference store st_{Ref} , including \mathfrak{S} and null . We define the concretization function γ of abstract heaps as:

$$(st_{\text{Ref}}, st_{\text{Val}}) \in \gamma(\overline{St}, \overline{con}, isSum) \Leftrightarrow \left(\begin{array}{l} \exists \alpha_{\text{Ref}} \in (references(st_{\text{Ref}}) \rightarrow \overline{\text{Ref}}). \\ \text{GraphEmbed}(\alpha_{\text{Ref}}, st_{\text{Ref}}, (\overline{St}, isSum)) \wedge \\ \text{ValueEmbed}(\alpha_{\text{Ref}}, (st_{\text{Ref}}, st_{\text{Val}}), \overline{con}) \end{array} \right)$$

That is, a concrete heap $(st_{\text{Ref}}, st_{\text{Val}})$ is in the concretization of an abstract heap $(\overline{St}, \overline{con}, isSum)$ iff there exists an embedding α_{Ref} (a function from concrete references to abstract nodes) such that the *shape* and the *values* of the concrete heap can be embedded into the abstract heap. These embeddings are expressed via the predicates *GraphEmbed* and *ValueEmbed*, which are defined as follows.

GraphEmbed holds if a given concrete reference store matches the shape of a given abstract heap, ignoring the value information. This is the case if \mathfrak{S} and null are the only concrete references that are abstracted to their abstract counterparts \mathfrak{S} and null (1), if, whenever multiple concrete references are abstracted to a single abstract reference, that abstract reference is a summary node (2), and if every concrete edge is represented by a (unique) abstract edge in the abstract heap (3):

$$\text{GraphEmbed}(\alpha_{\text{Ref}}, st_{\text{Ref}}, (\overline{St}, isSum)) \Leftrightarrow$$

$$\alpha_{\text{Ref}}^{-1}(\mathfrak{S}) = \{\mathfrak{S}\} \wedge \alpha_{\text{Ref}}^{-1}(\text{null}) = \{\text{null}\} \wedge \quad (1)$$

$$(\forall \overline{n} \in \text{img}(\alpha_{\text{Ref}}) \cdot |\alpha_{\text{Ref}}^{-1}(\overline{n})| > 1 \Rightarrow isSum(\overline{n})) \wedge \quad (2)$$

$$\forall r_1 \xrightarrow{f_r} r_2 \in st_{\text{Ref}} \cdot \alpha_{\text{Ref}}(r_1) \xrightarrow{f_r} \alpha_{\text{Ref}}(r_2) \in \overline{St} \quad (3)$$

where α_{Ref}^{-1} is the *preimage* of α_{Ref} (that is, it yields the set of concrete references abstracted by a given abstract reference).

ValueEmbed expresses that, for a given concrete reference store st_{Ref} , the value store st_{Val} matches all relevant abstract conditions in the abstract heap. Here, an abstract condition is *relevant* if it is associated with an abstract edge that corresponds to a concrete edge in st_{Ref} . In the definition below, we relate the concrete value store st_{Val} to each relevant abstract condition via a map s from abstract identifiers to sets of concrete values. For each concrete edge $r_1 \xrightarrow{f_r} r_2$ in the concrete reference store st_{Ref} , there is a map s in the concretization of the abstract condition of the corresponding abstract edge (4). The map s may constrain a concrete location (r, f_v) in three ways: via the abstract identifier $(\alpha_{\text{Ref}}(r), f_v)$, via the edge-local identifier (Src, f_v) if r is the source of the concrete edge, that is, $r = r_1$, and via the edge-local identifier (Trg, f_v) if r is the target of the concrete edge, that is, $r = r_2$. In all three cases, the map s must yield a set that contains the value v stored in the concrete value store for (r, f_v) (5). Finally, any concrete value store matches the relevant abstract conditions only if the conditions do not contradict each other, even on abstract locations that are not included in a given concrete heap. To ensure there are no such contradictions, s must be in the concretization of *all* relevant conditions, ignoring edge-local identifiers, which may denote different locations for different abstract edges. We use the operator \downarrow_{Loc} to project to the identifiers in Loc , that is, to remove edge-local identifiers (6).

$$\text{ValueEmbed}(\alpha_{\text{Ref}}, (st_{\text{Ref}}, st_{\text{Val}}), \overline{con}) \Leftrightarrow \forall r_1 \xrightarrow{f_r} r_2 \in st_{\text{Ref}} \cdot \exists s \in \gamma_{\overline{V}}(\overline{con}(\alpha_{\text{Ref}}(r_1) \xrightarrow{f_r} \alpha_{\text{Ref}}(r_2))) \cdot \quad (4)$$

$$\forall ((r, f_v) \mapsto v) \in st_{\text{Val}} \cdot \left(\begin{array}{l} v \in s(\alpha_{\text{Ref}}(r), f_v) \wedge \\ r = r_1 \Rightarrow v \in s(\text{Src}, f_v) \wedge \\ r = r_2 \Rightarrow v \in s(\text{Trg}, f_v) \end{array} \right) \wedge \quad (5)$$

$$s \downarrow_{\text{Loc}} \in \gamma_{\overline{V}} \left(\bigcap_{r'_1 \xrightarrow{f'_r} r'_2 \in st_{\text{Ref}}} \left(\overline{con}(\alpha_{\text{Ref}}(r'_1) \xrightarrow{f'_r} \alpha_{\text{Ref}}(r'_2)) \downarrow_{\text{Loc}} \right) \right) \quad (6)$$

Example. Fig. 5 shows the reference and value stores of two concrete heaps. The heap of the left) is in the concretization of the abstract heap in Fig. 2. For the embedding $\alpha_{\text{Ref}} = [\mathfrak{I} \mapsto \mathfrak{I}, \text{null} \mapsto \text{null}, r_1 \mapsto \bar{n}_0, r_2 \mapsto \bar{n}_0]$, *GraphEmbed* holds since \bar{n}_0 is a summary node and all three concrete edges have corresponding abstract edges. *ValueEmbed* also holds since the concrete value store satisfies the three relevant abstract conditions, and these conditions do not contradict each other.

In contrast, the heap on the right is *not* in the concretization of the abstract heap in Fig. 2. The graph embedding forces the embedding to be $\alpha_{\text{Ref}} = [\mathfrak{I} \mapsto \mathfrak{I}, \text{null} \mapsto \text{null}, r_1 \mapsto \bar{n}_0]$. Therefore, both edge-local identifiers (Src, val) and (Trg, val) on the abstract edge from \bar{n}_0 to \bar{n}_0 correspond to (r_1, val) , such that there is no value for (r_1, val) that satisfies the constraint $(\text{Trg}, \text{val}) = (\text{Src}, \text{val}) + 1$.

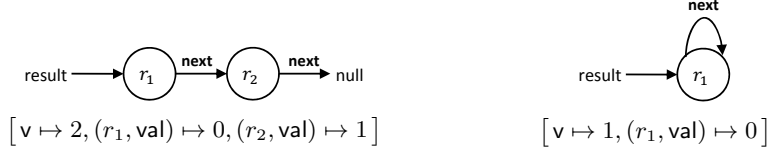


Fig. 5: Concrete heaps, consisting of a reference store, displayed on top, and a value store, displayed underneath. The heap on the left is in the concretization of the abstract heap in Fig. 2, whereas the heap on the right is not because it violates the condition that list nodes store increasing values.

In other words, any map s in the concretization of this constraint assigns different values to these edge-local identifiers and, thus, does not satisfy condition (5).

3.3 Join

The *join* operator $\sqcup_{\overline{\Sigma}}$ first computes an abstract reference store for the joined heaps and then the abstract conditions for the edges in this store.

Abstract Reference Store. An abstract heap can be viewed as a directed graph in which vertices are labeled as \mathfrak{S} , null , definite node other than \mathfrak{S} and null , or summary node; edges are labeled with reference fields. The vertex labels are used to avoid matching nodes in two heaps that cannot correspond (for instance, a summary node and a definite node). A *labeled heap graph* is a triple $g = (V, E, \eta) \in \mathbf{Graph}$, where $V \subseteq \overline{\mathbf{Ref}}$ is a set of vertices, $E \subseteq V \times \mathbf{Field}_{\mathbf{Ref}} \times V$ is a set of edges labeled with a reference field, and $\eta : V \rightarrow \{\mathfrak{S}, \text{Null}, \text{Def}, \text{Sum}\}$ is a labeling function on vertices. We assume a *strict total order* $<_{\mathbf{G}}$ on graphs that ensures in particular that $g_1 <_{\mathbf{G}} g_2$ if g_1 has fewer vertices than g_2 or the same number of vertices but fewer edges.

To improve performance, we define the join of two abstract reference stores such that it minimizes the size of the resulting store. Its structure is the minimum common supergraph of the two joined stores. Let g_1 and g_2 be graphs. Graph g is a *common supergraph* of g_1 and g_2 iff g_1 and g_2 are subgraph isomorphic to g with the isomorphisms \mathcal{I}_1 and \mathcal{I}_2 , respectively. We call g the *minimum common supergraph* (MCS) of g_1 and g_2 if there exists no other common supergraph that is smaller in the ordering $<_{\mathbf{G}}$. The procedure $\text{MCS}(g_1, g_2)$ yields the (unique) minimum common supergraph g of g_1 and g_2 as well as the corresponding subgraph isomorphisms \mathcal{I}_1 and \mathcal{I}_2 between g_1 and g , and g_2 and g , respectively. The problem of computing MCS can be reduced to the well-studied problem of finding the maximum common subgraph [3]. See Appendix A for the definitions of graph/subgraph isomorphism and maximum common subgraph. Intuitively, we can compute $(g, \mathcal{I}_1, \mathcal{I}_2) = \text{MCS}(g_1, g_2)$ by “gluing” to the maximum common subgraph of g_1 and g_2 those parts of g_1 and g_2 that are not in their maximum common subgraph.

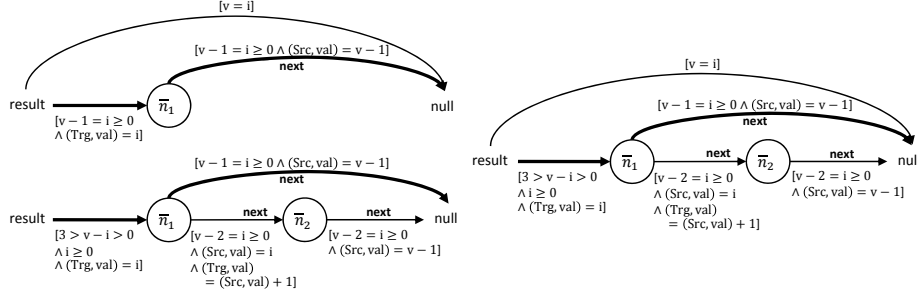


Fig. 6: The abstract heap graphs on the left occur before and after the second iteration of the fixed-point computation of the loop invariant (line 4) in Fig. 1. Joining them results in the heap on the right. Bold arrows indicate edges of the maximum common subgraph.

Let $\Pi : \overline{\text{Store}}_{\text{Ref}} \times \text{IsSummary} \rightarrow \text{Graph}$ be a *bijective* function from abstract stores to heap graphs. The abstract store and the `IsSummary` function of the join of $(\bar{St}_1, \bar{con}_1, isSum_1)$ and $(\bar{St}_2, \bar{con}_2, isSum_2)$ are $(\bar{St}, isSum) = \Pi^{-1}(MCS(\Pi(\bar{St}_1, isSum_1), \Pi(\bar{St}_2, isSum_2)) \downarrow_1)$, where \downarrow_1 denotes projection of a tuple on the first component, that is, the graph returned by *MCS*. Note that $\Pi(\bar{St}, isSum)$ includes both $\Pi(\bar{St}_1, isSum_1)$ and $\Pi(\bar{St}_2, isSum_2)$. Hence, the abstract reference store $(\bar{St}, isSum)$ subsumes the abstract reference stores $(\bar{St}_1, isSum_1)$ and $(\bar{St}_2, isSum_2)$.

Example. Fig. 6 shows on the left the abstract heap graphs g_1 and g_2 before and after the second iteration of the fixed-point computation of the loop invariant (line 4) in Fig. 1, and their join g on the right. Besides the special nodes \mathfrak{S} and `null`, the maximum common subgraph includes node \bar{n}_1 as well as the edges from `result` to \bar{n}_1 and from \bar{n}_1 to `null`. To this common subgraph, we add the remainder of g_1 (the edge from `result` to `null`) and the remainder of g_2 (\bar{n}_2 with its edges). Note that both g_1 and g_2 are subgraph isomorphic to g , where the isomorphism is the identity function.

Abstract Conditions. Consider an edge in the abstract store resulting from the join of two abstract heaps. We determine its abstract condition as follows. If the edge is in the maximum common subgraph of the joined heap graphs, its abstract condition is the join of the abstract conditions in the two heaps. Otherwise, the condition is the same as in the heap that contributed the edge, after applying the subgraph isomorphism.

As explained above, computing the minimum common supergraph $(g, \mathcal{I}_1, \mathcal{I}_2) = MCS(\Pi(\bar{St}_1, isSum_1), \Pi(\bar{St}_2, isSum_2))$ yields the subgraph isomorphisms \mathcal{I}_1 and \mathcal{I}_2 from $\Pi(\bar{St}_1, isSum_1)$ to g and from $\Pi(\bar{St}_2, isSum_2)$ to g , respectively. We define the function $rename_{ISO} : (\text{Ref} \rightarrow \text{Ref}) \times \bar{V} \rightarrow \bar{V}$ to rename the identifiers in $\overline{\text{Loc}}$ of a given abstract value state according to an isomorphism. Using this renaming, we define the join operator $\sqcup_{\bar{\Sigma}} : \bar{\Sigma} \times \bar{\Sigma} \rightarrow \bar{\Sigma}$ as

$$(\bar{St}_1, \bar{con}_1, isSum_1) \sqcup_{\bar{\Sigma}} (\bar{St}_2, \bar{con}_2, isSum_2) = (\bar{St}, \bar{con}, isSum)$$

where:

$$\begin{aligned}
(g, \mathcal{I}_1, \mathcal{I}_2) &= MCS(\Pi(\overline{St}_1, isSum_1), \Pi(\overline{St}_2, isSum_2)) \wedge \\
(\overline{St}, isSum) &= \Pi^{-1}(g) \wedge \\
\overline{con} &= \left[\overline{e} \mapsto \overline{s} \mid \overline{e} \in \overline{St} \wedge \overline{s} = \bigsqcup \left\{ \overline{s}' \mid \begin{array}{l} \exists i \in \{1, 2\} \cdot \exists (\overline{n}_1, f_r, \overline{n}_2) \in \overline{St}_i \cdot \\ \overline{e} = (\mathcal{I}_i(\overline{n}_1), f_r, \mathcal{I}_i(\overline{n}_2)) \wedge \\ \overline{s}' = rename_{ISO}(\mathcal{I}_i, \overline{con}_i(\overline{n}_1, f_r, \overline{n}_2)) \end{array} \right\} \right]
\end{aligned}$$

Computing the maximum common subgraph is NP-complete; however, most code fragments change only small portions of the abstract heap. Our implementation exploits this fact to compute the isomorphisms incrementally, usually in linear time.

Example. Consider the edge from **result** to \overline{n}_1 in the heap on the right of Fig. 6, which is in the maximum common subgraph of the heaps on the left. Hence, its abstract condition is the join of the conditions for those heaps (assuming a relational numerical domain). Since the constraint $v-1 = i$ in the top left abstract heap implies the constraint $3 > v-i > 0$ in the lower heap, the latter constraint is tracked by the result of the join operation; the other constraints of the joined conditions are identical and, thus, carried over to the result. Conversely, the edges from **result** to **null**, from \overline{n}_1 to \overline{n}_2 , and from \overline{n}_2 to **null** are not in the maximum common subgraph; their conditions come from the heap contributing the edges.

3.4 Widening

The above join operator does not guarantee the convergence of the analysis. In fact, the size of the abstract heap may grow at each application of join, and the abstract conditions may not stabilize. Therefore, we define a *widening* operator $\nabla_{\overline{\Sigma}} : \overline{\Sigma} \times \overline{\Sigma} \rightarrow \overline{\Sigma}$ that guarantees that the analysis reaches a fixed point in finite time (that is, terminates). In order to do so, the widening operator must bound the size of the abstract heap, which means that it has to merge nodes into summary nodes. This merging is controlled via a finite set of field access expressions \mathcal{W} , which is a parameter of the analysis and denotes references that the analysis should track separately. By default, \mathcal{W} is the set of local reference variables, but it can be extended to any set of field access expressions if desired. For all examples in our evaluation (Sec. 6), the analysis uses the default.

We perform widening in two steps. First, in both input heaps, we merge nodes that (i) are denoted by the same set of field access expressions from \mathcal{W} , and (ii) are reachable (via some access path) from the same set of local variables. Second, if the two heaps are isomorphic, we apply edge-wise widening to the abstract conditions; otherwise, we join them.

Merging Nodes. Our widening operator is parametric in a finite set of access expressions $\mathcal{W} \subseteq \mathbf{rAE}$, which denote objects that the user wants the analysis to track separately. Using access expressions for obtaining a bounded heap representation is standard [17,19,23]. Different heuristics and data-flow analyses can

be used to automatically determine \mathcal{W} . For instance, for all examples in our evaluation (Sec. 6), we used the set of local reference variables.

Based on this set of access expressions, we partition the nodes of an abstract heap into equivalence classes such that two nodes are in the same equivalence class iff they are reached via the same set of access expressions from \mathcal{W} and they are reachable from the same set of local variables. The latter condition improves precision, for instance, by preventing the widening from merging the tails of two separate lists into a single summary node. We define special equivalence classes for \mathfrak{S} and **null**.

Let $pointedBy_{\mathcal{W}} : (\overline{\text{Ref}} \times \overline{\text{Store}}_{\overline{\text{Ref}}}) \rightarrow \mathcal{P}(\mathcal{W})$ return the set of access expressions in \mathcal{W} leading to a given node in a given abstract reference store. Let $reach : \overline{\text{Ref}} \times \overline{\text{Store}}_{\overline{\text{Ref}}} \rightarrow \mathcal{P}(\text{Field}_{\overline{\text{Ref}}})$ be a function that returns the set of reference fields representing local variables from which a give node is reachable in a given abstract store. Formally:

$$reach(\overline{n}, \overline{St}) = \left\{ f \mid \begin{array}{l} \exists ae \in \mathbf{AE}_{\mathfrak{S}}, \langle \overline{e}_1, \dots, \overline{e}_n \rangle \in paths(ae, \overline{St}) . \\ target(\overline{e}_n) = \overline{n} \wedge field(e)_1 = f \end{array} \right\}$$

Let $partition : \overline{\text{Store}}_{\overline{\text{Ref}}} \rightarrow \mathcal{P}(\mathcal{P}(\overline{\text{Ref}}))$ be the a function that partitions the nodes of a given abstract store into equivalence classes as described above. The partitioning satisfies the following properties: (1) the partitions cover all nodes in the abstract store and there are separate partitions for \mathfrak{S} and **null**, (2) the partitions are non-empty and do not overlap, (3) all and only the nodes (other than \mathfrak{S} and **null**) nodes pointed by the same access expressions from \mathcal{W} and reachable from the same local variables are in the same partition. Note, that the above conditions define a *unique* partitioning and the number of partitions is bounded by the size of \mathcal{W} and the number of local variables. Formally: $partition(\overline{St}) = \{\overline{N}_1, \dots, \overline{N}_k\}$ where

$$\bigcup_{1 \leq i \leq k} \overline{N}_i = heapNodes(\overline{St}) \wedge \overline{N}_1 = \{\mathfrak{S}\} \wedge \overline{N}_2 = \{\text{null}\} \wedge \quad (1)$$

$$(\forall 1 \leq i, j \leq k . \overline{N}_i \neq \emptyset \wedge (i \neq j \Rightarrow \overline{N}_i \cap \overline{N}_j = \emptyset)) \wedge \quad (2)$$

$$\forall 3 \leq i, j \leq k, \overline{n}_1 \in \overline{N}_i, \overline{n}_2 \in \overline{N}_j .$$

$$i = j \iff \left(\begin{array}{l} pointedBy_{\mathcal{W}}(\overline{n}_1, \overline{St}) = pointedBy_{\mathcal{W}}(\overline{n}_2, \overline{St}) \wedge \\ reach(\overline{n}_1, \overline{St}) = reach(\overline{n}_2, \overline{St}) \end{array} \right) \quad (3)$$

After partitioning the nodes of an abstract heap, we merge all nodes in the same partition into a single (if necessary, summary) node. The merging of nodes may also replace several edges by a single edge from or to the merged node. The abstract condition for this edge is the join of the conditions of the edges it replaces.

We merge the nodes in the same partition into a single node, namely the smallest node of the partition, guaranteeing that the nodes from two different partitions are not merged into the same node. Let $min : \mathcal{P}(\overline{\text{Ref}}) \rightarrow \overline{\text{Ref}}$ return the smallest node from a given set of nodes. We also *fold* abstract identifiers in the abstract conditions of the heap corresponding to the merged nodes. Let

$getFolds : \overline{\text{Store}}_{\text{Ref}} \rightarrow \mathcal{P}(\mathcal{P}(\overline{\text{Loc}}) \times \overline{\text{Loc}})$ be a function that describes how the abstract identifiers should be folded. Recall, that $fold(R, \bar{s})$ returns a value state in which for any $(\bar{I}, \bar{i}) \in R$, identifiers in \bar{I} are replaced with \bar{i} and the value of \bar{i} is the least upper-bound of the values of identifiers in \bar{I} . Note, that we do not need to take care of the edge-local identifiers. Formally:

$$getFolds(\overline{St}) = \{(R, (min(\overline{N}_i), f_v)) \mid \{\overline{N}_1, \dots, \overline{N}_k\} = partition(\overline{St}) \wedge f_v \in \text{Field}_{\text{Val}} \wedge 1 \leq i \leq k \wedge R = \{(\overline{n}_i, f_v) \in \overline{\text{Loc}} \mid \overline{n}_i \in \overline{N}_i\}\}$$

Now, we are in the position to define a function $merge_{\overline{\Sigma}} : \overline{\Sigma} \rightarrow \overline{\Sigma}$ that formalizes what it means to merge nodes of an abstract heap with respect to given \mathcal{W} . First, we compute the map ϕ that describes how the nodes are merged, and subsequently, we merge the nodes in the abstract store. Then, we compute which abstract identifiers of abstract conditions should be merged and merge them. Merging nodes may also lead to merging of several edges. Therefore, we need to join the conditions of merged edges. Finally, we update the map defines nodes as summary or definite. If the partition contained only a single definite node, the node stays definite. If we have merged several nodes or the partition already contained a summary node, the node into which we have merged the partition is a summary node. Formally: $merge_{\overline{\Sigma}}(\overline{St}, \overline{con}, isSum) = (\overline{St}', \overline{con}', isSum')$ where

$$\begin{aligned} \phi &= \left[\overline{n} \mapsto min(\overline{N}_i) \mid \left\{ \overline{N}_1, \dots, \overline{N}_k \right\} = partition(\overline{St}) \wedge \right. \\ &\quad \left. \exists 1 \leq i \leq k . \overline{n} \in \overline{N}_i \right] \wedge \\ \overline{St}' &= \{(\phi(\overline{n}_1), f_r, \phi(\overline{n}_2)) \mid (\overline{n}_1, f_r, \overline{n}_2) \in \overline{St}\} \wedge \\ \overline{con}' &= [(\overline{n}'_1, f_r, \overline{n}'_2) \mapsto \overline{s} \mid (\overline{n}'_1, f_r, \overline{n}'_2) \in \overline{St}' \wedge \\ &\quad \overline{s} = \bigsqcup \{fold(getFolds(\overline{St}), \overline{con}(\overline{n}_1, f_r, \overline{n}_2)) \mid \\ &\quad (\overline{n}_1, f_r, \overline{n}_2) \in \overline{St} \wedge \phi(\overline{n}_1) = \overline{n}'_1 \wedge \phi(\overline{n}_2) = \overline{n}'_2\}] \wedge \\ isSum' &= \left[\overline{n} \mapsto \neg \left(\exists 1 \leq i \leq k . \{\overline{n}\} = \overline{N}_i \right) \mid \overline{n} \in heapNodes(\overline{St}') \wedge \right. \\ &\quad \left. \wedge \neg isSum(\overline{n}) \right] \mid \{\overline{N}_1, \dots, \overline{N}_k\} = partition(\overline{St}) \end{aligned}$$

From Merge to Widening. Once we merged nodes in the heaps to which we want to apply widening, we need to perform edge-wise widening of the abstract conditions of the merged heaps. The widening of the conditions on the edges of the heaps can be performed only if the merged heaps are structurally the same (isomorphic), otherwise there are edges in one heap that do not correspond to any edge in the other heap, hence, their conditions can not be widened. Therefore, we need to be able to realize whether two merged heaps are isomorphic.

We observe that if \mathcal{I} is an isomorphism from a heap $(\overline{St}_2, \overline{con}_2, isSum_2)$ to a heap $(\overline{St}_1, \overline{con}_1, isSum_1)$, then for any $\overline{n} \in dom(\mathcal{I})$, the nodes \overline{n} and $\mathcal{I}(\overline{n})$ must be reachable from the same set of local variables and the same set of access expressions in \mathcal{W} must lead to them in their respective abstract stores. Furthermore, it must be the case that $isSum_2(\overline{n}) = isSum_1(\mathcal{I}(\overline{n}))$ and that if $\overline{n} \in \{\mathfrak{S}, \text{null}\}$ then $\overline{n} = \mathcal{I}(\overline{n})$. Due to the uniqueness of the partitioning of nodes of an abstract heap, any two different nodes of the abstract heap $merge_{\overline{\Sigma}}(\overline{\sigma})$

other than \mathfrak{S} and null are reachable from different sets of local variables, or different sets of access expressions in \mathcal{W} lead to them. Therefore, we can find a candidate for an isomorphism between two heaps in which we have merged the nodes just by looking at the set of access expressions from \mathcal{W} leading to the node of the heaps and local variables from which the nodes are reachable. Let $find_{\text{ISO}} : ((\overline{\text{Store}}_{\overline{\text{Ref}}} \times \text{IsSummary}) \times (\overline{\text{Store}}_{\overline{\text{Ref}}} \times \text{IsSummary})) \rightarrow (\overline{\text{Ref}} \rightarrow \overline{\text{Ref}})$ return a possible witness (candidate) for the isomorphism from the second abstract store to the first. Formally:

$$find_{\text{ISO}}((\overline{St}_1, isSum_1), (\overline{St}_2, isSum_2)) = \left[\begin{array}{l} \overline{n}_1 \in heapNodes(\overline{St}_1) \wedge \overline{n}_2 \in heapNodes(\overline{St}_2) \wedge \\ isSum_1(\overline{n}_1) = isSum_2(\overline{n}_2) \wedge (\overline{n}_2 \in \{\mathfrak{S}, \text{null}\} \Rightarrow \overline{n}_2 = \overline{n}_1) \wedge \\ pointedBy_{\mathcal{W}}(\overline{n}_1, \overline{St}) = pointedBy_{\mathcal{W}}(\overline{n}_2, \overline{St}) \wedge \\ reach(\overline{n}_1, \overline{St}) = reach(\overline{n}_2, \overline{St}) \end{array} \right]$$

In contrast with NP-complete problem of searching for isomorphisms between two arbitrary graphs, the complexity of $find_{\text{ISO}}$ is linear in the number of nodes in given abstract stores and this number is bound by the size of \mathcal{W} and the number of local variables. Hence, once we fix \mathcal{W} and the set of local variables, $find_{\text{ISO}}$ is of constant complexity. Furthermore, if $find_{\text{ISO}}((\overline{St}_1, isSum_1), (\overline{St}_2, isSum_2))$ is not the isomorphism from \overline{St}_2 to \overline{St}_1 (i.e. not every node from one heap is mapped to a node in the other heap), then there exists no such isomorphism.

Even if the abstract stores of the merged heaps are structurally the same, we might not be able to perform widening on the abstract conditions of the corresponding edges, as the widening of value states is defined only for states that contain the same set of abstract identifiers. Hence, we define a function $compatible : (\overline{\Sigma} \times \overline{\Sigma}) \rightarrow \{\text{true}, \text{false}\}$ that checks whether the merged heaps are isomorphic and whether the abstract conditions of corresponding edges contain the same abstract identifiers with respect to the isomorphism. Formally:

$$\begin{aligned} compatible((\overline{St}_1, \overline{con}_1, isSum_1), (\overline{St}_2, \overline{con}_2, isSum_2)) &= \exists \mathcal{I} : \overline{\text{Ref}} \rightarrow \overline{\text{Ref}} . \\ \mathcal{I} &= find_{\text{ISO}}((\overline{St}_1, isSum_1), (\overline{St}_2, isSum_2)) \wedge dom(\mathcal{I}) = heapNodes(\overline{St}_2) \wedge \\ img(\mathcal{I}) &= heapNodes(\overline{St}_1) \wedge \\ (\forall \overline{e}_2 = (\overline{n}_2, f, \overline{n}'_2) \in \overline{St}_2 . \exists \overline{e}_1 = (\mathcal{I}(\overline{n}_2), f, \mathcal{I}(\overline{n}'_2)) \in \overline{St}_1 . \\ &\quad get_{\overline{\text{ID}}}(\overline{con}(\overline{e}_1)) = get_{\overline{\text{ID}}}(\text{rename}_{\text{ISO}}(\mathcal{I}, \overline{con}(\overline{e}_2)))) \wedge \\ (\forall \overline{e}_1 = (\overline{n}_1, f, \overline{n}'_1) \in \overline{St}_1 . \exists \overline{e}_2 = (\mathcal{I}^{-1}(\overline{n}_1), f, \mathcal{I}^{-1}(\overline{n}'_1)) \in \overline{St}_2 . \\ &\quad get_{\overline{\text{ID}}}(\overline{con}(\overline{e}_1)) = get_{\overline{\text{ID}}}(\text{rename}_{\text{ISO}}(\mathcal{I}, \overline{con}(\overline{e}_2)))) \end{aligned}$$

where $get_{\overline{\text{ID}}} : \overline{V} \rightarrow \mathcal{P}(\overline{\text{ID}})$ returns the set of abstract identifiers of a given abstract value state.

If the merged heaps are compatible, we can apply edge-wise widening to the conditions of the heaps. Let $widen_{\overline{\text{Cond}}} : (\overline{\text{Cond}} \times \overline{\text{Cond}}) \rightarrow \overline{\text{Cond}}$ be a function that applies widening edge-wise to conditions of compatible heaps. Formally,

$$widen_{\overline{\text{Cond}}}(\overline{con}_1, \overline{con}_2) = [\overline{e} \mapsto \overline{con}_1(\overline{e}) \nabla \overline{con}_2(\overline{e}) \mid \overline{e} \in dom(\overline{con}_1) \cap dom(\overline{con}_2)]$$

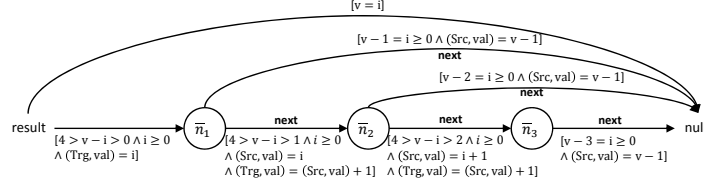


Fig. 7: Heap before the fourth iteration of the fixed-point computation of the loop invariant in Fig. 1.

Now, we are in the position to define the widening operator $\nabla_{\bar{\Sigma}} : (\bar{\Sigma} \times \bar{\Sigma}) \rightarrow \bar{\Sigma}$. First, we compute an intermediate heap that is the join of the heaps we want to widen. We do this in order to guarantee that all the edges from both heaps are considered when we subsequent merge of nodes is performed. After merging the nodes in the first heap and the intermediate heap, we check whether these heaps are compatible for widening of their abstract conditions. This means that the heaps are isomorphic and the identifiers of the abstract conditions are stabilized. If they are, we rename identifiers in the conditions of the intermediate merged heap according to the isomorphism \mathcal{I} that is the witness to the compatibility of the merged heaps, and we widen the conditions on the edges. On the other hand, if the merged heaps are not compatible, we return the intermediate heap as it is as upper bound on the heaps we want to widen. Formally,

$$(\bar{\sigma}_1) \nabla_{\bar{\Sigma}} (\bar{\sigma}_2) = \begin{cases} (\bar{st}'_1, \text{widen}_{\text{Cond}}(\bar{con}'_1, \bar{con}'_2), \text{isSum}'_1) & \text{if } \text{compatible}(\bar{\sigma}'_1, \bar{\sigma}'_2) \\ \bar{\sigma}'_2 & \text{otherwise} \end{cases}$$

where $\bar{\sigma}'_1 = (\bar{st}'_1, \bar{con}'_1, \text{isSum}'_1) = \text{merge}_{\bar{\Sigma}}(\bar{\sigma}_1)$, $\bar{\sigma}'_2 = (\bar{st}'_2, \bar{con}'_2, \text{isSum}'_2) = \text{merge}_{\bar{\Sigma}}(\bar{\sigma}_1 \sqcup \bar{\sigma}_2)$, in case of the merged heaps being compatible, \mathcal{I} is the witness and $\bar{con}'_2 = [(\mathcal{I}(\bar{n}_1), f_r, \mathcal{I}(\bar{n}_2)) \mapsto \text{rename}_{\text{ISO}}(\mathcal{I}, \bar{con}'_2(\bar{n}_1, f_r, \bar{n}_2)) \mid (\bar{n}_1, f_r, \bar{n}_2) \in \bar{st}'_2]$.

Termination. Since \mathcal{W} is finite and there are only finitely many local variables, there are only finitely many equivalence classes for the merging of nodes. Therefore, the abstract heap resulting from merging has a bounded number of nodes. When the two heaps after merging of nodes are isomorphic, termination of widening follows from the fact that there are only finitely many edges and that widening of the value analysis is assumed to terminate. When they are not isomorphic, the heaps are joined. However, the join can be applied only finitely many times, since there are only finitely many non-isomorphic structures that can be constructed from a fixed number of nodes, and after every application of join, the size of the abstract structure grows, as join is a monotone operator. Furthermore, also the set of identifiers tracked by conditions on edges grows with every application of join to conditions and there are only finitely many identifiers (since there are only finitely many nodes in the heap and $\text{Field}_{\text{Val}}$ is finite).

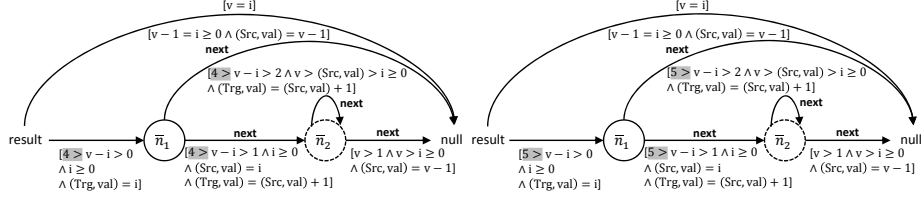


Fig. 8: Heaps with merged nodes before the fourth (left) and fifth (right) iteration of the fixed-point computation. The heaps differ only in the highlighted constraints.

Example. Suppose we widen the abstract heap before the fourth iteration of the fixed-point computation for the loop in Fig. 1 with the heap before the fifth iteration. The abstract heap before the fourth iteration is displayed in Fig. 7; the heap before the fifth iteration looks similar, but has four definite nodes.

In the first step, widening merges nodes using the default $\mathcal{W} = \{\text{result}\}$. In the heap from Fig. 7, we merge \bar{n}_2 and \bar{n}_3 into a single summary node \bar{n}_2 since they are (i) denoted by the same set of field access expressions from \mathcal{W} (the empty set since **result** denotes neither \bar{n}_2 nor \bar{n}_3), and (ii) they are reachable from the same set of local variables ($\{\text{result}\}$). However, \bar{n}_1 is denoted by a different set of field access expressions ($\{\text{result}\}$), and therefore not merged with \bar{n}_2 and \bar{n}_3 . The edges from \bar{n}_2 and \bar{n}_3 to **null** are also merged, and their conditions are joined. The resulting heap is shown on the left of Fig. 8. Merging nodes in the heap before the fifth iteration (not shown) results in the heap on the right of Fig. 8. Note that these heaps are isomorphic, that is, the heap shape has stabilized.

In the second step, since the heaps after merging are isomorphic, we apply edge-wise widening to the abstract conditions. This step removes the upper bound on $v - i$, but leaves all other constraints unaffected, that is, the abstract conditions have stabilized. The resulting heap is shown in Fig. 4; it represents the loop invariant of the program from Fig. 1.

4 Evaluation of Expressions

4.1 Evaluation of Reference Expressions

Here, we formalize the evaluation of reference expressions (**rexp**). A reference expression is either **null**, an allocation of a new object (**new C**), or a reference field access expression (**rAE**). We first formalize the evaluation of **new C** and the evaluation of **rAE**. Finally, we define the evaluation function of reference expressions.

We assume that the underlying value analysis \bar{V} provides the following operators based on well-known definitions [15]: (i) *add* : $\mathcal{P}(\text{ID}) \times \bar{V} \rightarrow \bar{V}$, which introduces a given set of identifiers with value \top into a given state, (ii) *drop* : $\mathcal{P}(\text{ID}) \times \bar{V} \rightarrow \bar{V}$, which projects out a given set of identifiers from a state, and (iii) *expand* : $\mathcal{P}(\text{ID} \times \text{ID}) \times \bar{V} \rightarrow \bar{V}$, which extends a given state such that for any (\bar{id}_1, \bar{id}_2) of a given set of pairs of identifiers, the abstract value of \bar{id}_2 is the

same as the abstract value of \overline{id}_1 , and the values of the identifiers present in the original state remain unchanged.

We partition the set of edge-local identifiers into *source* edge-local identifiers $\overline{\text{Eld}}_{\text{Src}} = \{\text{Src}\} \times \text{Field}_{\text{Val}}$ and *target* edge-local identifiers $\overline{\text{Eld}}_{\text{Trg}} = \{\text{Trg}\} \times \text{Field}_{\text{Val}}$.

In the evaluation of both reference and value expressions, we will have to determine a value state that approximates the value stores of all heaps represented by a given abstract heap. We call such a value state the *general condition* of an abstract heap. It is defined as follows. *Each* abstract condition associated with an edge of an abstract heap approximates the value store of *all* concrete heaps in which the edge exists. Therefore, it is sufficient to consider any abstract reference \bar{n} that abstracts a concrete reference that exists in every concrete heap (such as \mathfrak{S}) and any field f . The join of the conditions of the edges $(\bar{n}, f, -)$ is a value state that approximates the value stores of all concrete stores (since $\bar{n}.f$ must have some value in each of them). In the definition below, we take the join over all fields for simplicity, and drop the edge-local identifiers since the general condition is not associated with any particular edge and, hence, the edge-local identifiers have no meaning. Formally, the general condition $general_{\bar{V}} : \bar{\Sigma} \rightarrow \bar{V}$ of an abstract heap is defined as $general_{\bar{V}}(\bar{St}, \bar{con}, -) = \bigsqcup_{\bar{e} \in \bar{St}} \{drop(\overline{\text{Eld}}, \bar{con}(\bar{e})) \mid \mathfrak{S} = source(\bar{e})\}$.

Object Allocation. An abstract heap $\bar{\sigma}'$ soundly approximates an allocation of a new object in abstract heap $\bar{\sigma}$ if and only if for each concrete heap $h \in \gamma(\bar{\sigma})$ and for each concrete heap h' resulting from the allocation of a new object in h , the concretization of $\bar{\sigma}'$ contains h' , that is, $h' \in \gamma(\bar{\sigma}')$. Let $alloc_{\bar{\Sigma}} : \bar{\Sigma} \rightarrow (\mathcal{P}((\overline{\text{Ref}} \times \bar{V})) \times \bar{\Sigma})$ be a function that for a given abstract heap returns (i) a singleton set³ containing the abstract reference representing the newly allocated object and the abstract condition approximating value stores of heaps resulting from the allocation, and (ii) an abstract heap which soundly approximates object allocation in the given heap. Formally, $alloc_{\bar{\Sigma}}(\bar{St}, \bar{con}, isSum) = (\{(\bar{n}, \bar{s})\}, (\bar{St}', \bar{con}', isSum'))$ where

$$\bar{n} \in (\overline{\text{Ref}} \setminus (\{\bar{n}' \mid (\bar{n}', -, -) \in \bar{St} \vee (-, -, \bar{n}') \in \bar{St}\} \cup \{\text{null}, \mathfrak{S}\})) \wedge \quad (1)$$

$$isSum' = isSum[\bar{n} \mapsto \text{false}] \wedge \quad (2)$$

$$\bar{I} = \{\bar{n}\} \times \text{Field}_{\text{Val}} \wedge \quad (3)$$

$$\bar{s} = add(\bar{I} \cup \overline{\text{Eld}}, general_{\bar{V}}(\bar{St}, \bar{con}, isSum)) \wedge \quad (4)$$

$$\bar{St}' = \bar{St} \cup \{(\bar{n}, f_r, \text{null}) \mid f_r \in \text{Field}_{\text{Ref}}\} \wedge \quad (5)$$

$$\bar{con}' = \left[\bar{e} \mapsto \begin{cases} add(\bar{I}, \bar{con}(e)) & \text{if } \bar{e} \in \bar{St} \\ drop(\overline{\text{Eld}}_{\text{Trg}}, \bar{s}) & \text{otherwise} \end{cases} \middle| \bar{e} \in \bar{St}' \right] \quad (6)$$

The newly allocated object is represented by a fresh non-null node \bar{n} (1), which is definite (2). The analysis builds up the value state \bar{s} approximating the value store of all concrete heaps in which the new object has been allocated (using the general condition defined above). Consequently, \bar{s} must contain the identifiers \bar{I}

³ Returning a singleton set rather than a single pair simplifies the definition of $\overline{eval}_{\text{rexp}}$ in Sec. 4.1.

for the value fields of \bar{n} (3). Since the allocation of an object can appear only on the right-hand side of an assignment, \bar{s} will participate in the condition of an edge that represents the assignment (see line (6) of the abstract semantics of reference assignment in Sec. 5.1). For this purpose, we also introduce edge-local identifiers into \bar{s} (4). We assume that the concrete semantics of object allocation assigns **null** to reference fields of the newly allocated object. Accordingly, the analysis updates the abstract store such that the reference fields of \bar{n} point to **null** (5). Finally, the abstract conditions are updated: for each edge that is already present, we add identifiers \bar{I} to their existing conditions. For every new edge, we add the condition \bar{s} in which the target edge-local identifiers are not present (since the target of the edge is **null**) (6).

Reference Field Access Expressions. Since the abstract store includes disjunctive information, a reference field access expression $\mathfrak{F}.f_1 \dots .f_n \in \text{rAE}$ may correspond to multiple paths in the abstract heap. When evaluating $\mathfrak{F}.f_1 \dots .f_n$, all these paths must be considered. Let $paths : (\text{rAE} \times \text{Store}_{\text{Ref}}) \rightarrow \mathcal{P}((\text{Ref} \times \text{Field}_{\text{Ref}} \times \text{Ref})^*)$ be a function that returns the set of all paths in a given abstract store for a given reference field access expression. Formally:

$$paths(\mathfrak{F}.f_1 \dots .f_n, \bar{St}) = \{ \langle \bar{e}_1, \dots, \bar{e}_n \rangle \mid \forall 1 \leq i \leq n \cdot \bar{e}_i \in \bar{St} \wedge field(\bar{e}_i) = f_i \wedge source(\bar{e}_1) = \mathfrak{F} \wedge ((1 < i) \Rightarrow (source(\bar{e}_i) = target(\bar{e}_{i-1}))) \}$$

Each edge in a path is associated with an abstract condition that approximates the value stores of all concrete heaps in which the edge exists. Therefore, the path is feasible only in those concrete heaps whose value stores are approximated by the abstract conditions associated with *all* edges of the path. However, since the abstract condition of an edge contains edge-local identifiers for the source and the target of this particular edge, it is not generally sound to simply take the greatest lower bound of all abstract conditions of the path. If \bar{e}_i and \bar{e}_{i+1} are two consecutive edges in the path then the target edge-local identifiers in the abstract condition of \bar{e}_i correspond to the source edge-local identifiers in the abstract condition of \bar{e}_{i+1} . To adjust the edge-local identifiers, we define a function $\text{TrgToSrc} : \bar{V} \rightarrow \bar{V}$, which converts the target edge-local identifiers of a given abstract condition to source edge-local identifiers and sets the target edge-local identifiers to \top . Formally:

$$\text{TrgToSrc}(\bar{s}) = add(\bar{\text{Eld}}_{\text{Trg}}, drop(\bar{\text{Eld}}_{\text{Trg}}, expand(R, drop(\bar{\text{Eld}}_{\text{Src}}, \bar{s}))))$$

where $R = \{((\text{Trg}, f_v), (\text{Src}, f_v)) \mid f_v \in \text{Field}_{\text{Val}}\}$. In the above definition, we convert the target edge-local identifiers to source local-identifiers by first removing the source edge-local identifiers from the value state (via *drop*), and then re-introducing them with the abstract values of the corresponding target edge-local identifiers (via *expand*). Since the source edge-local identifiers assumed the role of the target edge-local identifiers, the target edge-local identifiers are removed from the state and then re-introduced with the value \top (via *add*).

We now define a function $pathCondition : (\bar{V} \times (\text{Ref} \times \text{Field}_{\text{Ref}} \times \text{Ref})^* \times \text{Cond}) \rightarrow \bar{V}$, which, for a given abstract condition and a given path, returns the

abstract condition under which the path is feasible and does not contradict the condition. Furthermore, in case the last node of the given path is non-null, the target edge-local identifiers of the resulting abstract condition refer to the value fields of the last node of the path. Formally:

$$\text{pathCondition}(\bar{s}, \langle \bar{e}_1, \dots, \bar{e}_n \rangle, \overline{con}) = \begin{cases} \text{add}(\overline{\text{Eld}}_{\text{src}}, \text{drop}(\overline{\text{Eld}}_{\text{src}}, \bar{s})) & \text{if } \langle \bar{e}_1, \dots, \bar{e}_n \rangle = \varepsilon \\ \text{pathCondition}(\text{Trg To Src}(\bar{s}) \sqcap \overline{con}(\bar{e}_1), \langle \bar{e}_2, \dots, \bar{e}_n \rangle, \overline{con}) & \text{otherwise} \end{cases}$$

where ε represents an empty path. In the above definition, \bar{s} represents the value information accumulated from the abstract conditions of already processed edges. In the last line of the definition, we obtain the next value of this accumulator by taking the greatest lower bound of the current \bar{s} (with appropriate adjustments of edge local identifiers) and the abstract condition of the next unprocessed edge.

Now we define an evaluation function for reference field access expressions $\overline{\text{eval}}_{\text{rAE}} : (\text{rAE} \times \bar{\Sigma}) \rightarrow \mathcal{P}(\text{Ref} \times \bar{V})$, which returns a set of pairs of abstract references and abstract value states. In each pair, the reference is the last node of a paths that may represent the given reference field access expressions, and the value state describes under which condition this path is feasible. In each pair, the target edge-local identifiers of the abstract condition refer to value fields of the abstract reference. Note that if the abstract condition is \perp , the path is not feasible and such a path may be discarded. Formally:

$$\overline{\text{eval}}_{\text{rAE}}(ap, (\bar{St}, \overline{con}, \text{isSum})) = \{(\bar{n}, \bar{s}) \mid \exists \langle \bar{e}_1, \dots, \bar{e}_n \rangle \in \text{paths}(ap, \bar{St}) \cdot \text{target}(\bar{e}_n) = \bar{n} \wedge \bar{s} = \text{pathCondition}(\top, \langle \bar{e}_1, \dots, \bar{e}_n \rangle, \overline{con}) \wedge \bar{s} \neq \perp\}$$

Reference Expressions. We define an evaluation function $\overline{\text{eval}}_{\text{rexp}} : ((\text{rexp} \cup \{\mathfrak{S}\}) \times \bar{\Sigma}) \rightarrow (\mathcal{P}(\text{Ref} \times \bar{V}) \times \bar{\Sigma})$ for reference expressions and also \mathfrak{S} , which simplifies the treatment of assignments to expressions of the form $\mathfrak{S}.f$. $\overline{\text{eval}}_{\text{rexp}}$ returns (i) the abstract references that a given reference expression (or \mathfrak{S}) may represent together with the abstract conditions under which it may represent these abstract references. In these conditions, the target edge-local identifiers must correspond to value fields of the resulting abstract references. And (ii) the abstract heap that results from evaluating the given expression, since allocation has side-effects. Formally:

$$\overline{\text{eval}}_{\text{rexp}}(re, \bar{\sigma}) = \begin{cases} \text{alloc}_{\bar{\Sigma}}(\bar{\sigma}) & \text{if } re = \text{new } C \\ \{(\{\text{null}, \text{add}(\overline{\text{Eld}}_{\text{src}}, \text{general}_{\bar{V}}(\bar{\sigma}))\}, \bar{\sigma}) & \text{if } re = \text{null} \\ (\overline{\text{eval}}_{\text{rAE}}(re, \bar{\sigma}), \bar{\sigma}) & \text{if } re \in \text{rAE} \\ \{(\{\mathfrak{S}, \text{add}(\overline{\text{Eld}}_{\text{src}}, \text{expand}(R, \text{general}_{\bar{V}}(\bar{\sigma})))\}, \bar{\sigma}) & \text{if } re = \mathfrak{S} \end{cases}$$

where $R = \{((\mathfrak{S}, f_v), (\text{Trg}, f_v)) \mid f_v \in \text{Field}_{\text{val}}\}$. Since the abstract conditions returned by $\overline{\text{eval}}_{\text{rexp}}$ may participate in conditions of edges (see line (6) of the abstract semantics of reference assignment in Sec. 5.1), we need to make sure that they include the edge-local identifiers.

4.2 Evaluation of Value Expressions

Here, we formalize the evaluation of value expressions (**vexp**). A value expression is either a literal (**n**), a value field access expression (**vAE**), or a binary expression (**vexp** **op** **vexp**). First, we formalize the evaluation of **vAE**. Then we define the evaluation function for value expressions.

Value Field Access Expressions. For a value field access expression $p.f_v$, the receiver expression p is either \exists or a reference field access expression. In the latter case, it may correspond to multiple paths in the abstract heap. Therefore, there may be multiple ways to evaluate $p.f_v$, and the evaluation function must consider each of them.

We use function $\overline{eval}_{\text{rexp}}$ (Sec. 4.1) to determine the set of abstract references representing p , and for each of them an abstract condition under which p represents this reference. The edge-local identifier (Trg, f_v) in these conditions refers to $p.f_v$, and, therefore, is important for the evaluation. However, in Sec. 5.2, we will detach these conditions from their edges such that edge-local identifiers lose their meaning. To preserve the information about $p.f_v$, we retain it using a separate, non-edge-local identifier before removing the edge-local identifiers. For this purpose, we extend the set of abstract identifiers $\overline{\text{ID}}$ of the value analysis with *value field access identifiers*, that is, $\overline{\text{ID}}_{\text{vAE}} = \overline{\text{ID}} \cup \text{vAE}$. These identifiers are used only temporarily for intermediate steps, but do not occur in any value states before or after a value assignment is performed.

Function $\overline{eval}_{\text{vAE}} : (\text{vAE} \times \overline{\Sigma}) \rightarrow \mathcal{P}(\overline{\text{V}})$ yields the set of value states corresponding to different ways a given value field access expression may be evaluated. Formally:

$$\overline{eval}_{\text{vAE}}(p.f_v, \overline{\sigma}) = \left\{ \text{drop}(\overline{\text{Eld}}, \overline{s}) \mid \begin{array}{l} \exists (\overline{n}, \overline{s}') \in \overline{eval}_{\text{rexp}}(p, \overline{\sigma}) \cdot \overline{n} \neq \text{null} \wedge \\ \overline{s} = \text{expand}(\{((\text{Trg}, f_v), p.f_v)\}, \overline{s}') \end{array} \right\}$$

For a given value field access expression $p.f_v$, the function first computes the abstract references representing p , together with the value states under which p represents them. Abstract null references are excluded since they would lead to null-pointer dereferencing when evaluating $p.f_v$ and, thus, abortion of the execution. Then, for each abstract condition, the function renames the identifier (Trg, f_v) to the new value field access identifier $p.f_v$ (via *expand*). Finally, since the resulting value state describes conditions on the entire heap, the function removes all edge-local identifiers.

Value Expressions. We define an evaluation function for value expressions $\overline{eval}_{\text{vexp}} : (\text{vexp} \times \overline{\Sigma}) \rightarrow \mathcal{P}(\overline{\text{V}})$, which returns the set of value states corresponding to every possible way a given value expression may be evaluated in a given abstract heap. Each resulting value state contains value field access identifiers for

all value field access expressions presents in the given value expression. Formally:

$$\overline{eval}_{\text{vexp}}(ve, \bar{\sigma}) = \begin{cases} \{general_{\bar{V}}(\bar{\sigma})\} & \text{if } ve \in \mathbf{n} \\ \overline{eval}_{\text{vAE}}(ve, \bar{\sigma}) & \text{if } ve \in \text{vAE} \\ \left\{ meet_{\text{vAE}}(\bar{s}_1, \bar{s}_2) \left| \begin{array}{l} \bar{s}_1 \in \overline{eval}_{\text{vexp}}(ve_1, \bar{\sigma}) \wedge \\ \bar{s}_2 \in \overline{eval}_{\text{vexp}}(ve_2, \bar{\sigma}) \end{array} \right. \right\} & \text{if } ve_1 \langle \text{op} \rangle ve_2 = ve \end{cases}$$

For literals \mathbf{n} , the function yields the general condition since a literal can be evaluated in any heap. For value field access expressions, the function applies $\overline{eval}_{\text{vAE}}$ from the previous section. The evaluation of binary expressions works as follows.

Recursively evaluating the expressions ve_1 and ve_2 yields sets of abstract states that represent all ways of evaluating these expressions. The resulting state for the binary expression is essentially the point-wise application of the greatest lower bound to each combination of states from these sets. However, these states may in general contain different value field access identifiers \bar{I}_1 and \bar{I}_2 . To avoid information loss when applying the greatest lower bound, we first enlarge all involved value states to include all identifiers in $\bar{I}_1 \cup \bar{I}_2$. Let $get_{\text{vAE}} : \bar{V} \rightarrow \mathcal{P}(\text{ID}_{\text{vAE}})$ yield a set of value field access identifiers of a given value state. Function $meet_{\text{vAE}} : \bar{V} \times \bar{V} \rightarrow \bar{V}$ yields the value state that is the greatest lower bound of the given value states and contains all value field access identifiers from these states:

$$meet_{\text{vAE}}(\bar{s}_1, \bar{s}_2) = add(get_{\text{vAE}}(\bar{s}_2) \setminus get_{\text{vAE}}(\bar{s}_1), \bar{s}_1) \sqcap add(get_{\text{vAE}}(\bar{s}_1) \setminus get_{\text{vAE}}(\bar{s}_2), \bar{s}_2)$$

5 Abstract Semantics

In this section, we formalize the semantics of reference and value assignments.

5.1 Reference Assignment

An abstract store includes disjunctive information. Therefore, for a reference assignment $\mathbf{p.f_r} = \mathbf{rhs}$, there may be several abstract references for the receiver \mathbf{p} and the right-hand side \mathbf{rhs} , which may be reached through different paths with different value conditions. The value states on the edges along each path specify the conditions under which \mathbf{p} and \mathbf{rhs} evaluate to a particular abstract reference. The abstract semantics for reference assignments adds an abstract edge for each possible combination of receiver \mathbf{p} and right-hand side \mathbf{rhs} , with an abstract condition that reflects when this combination exists.

The rule below formalizes reference assignments of the form $\mathbf{p.f_r} = \mathbf{rhs}$, where $\mathbf{p.f_r} \in \text{rAE}$ and $\mathbf{rhs} \in \text{rexp}$. Since we encode local variables as fields of a special reference \mathfrak{S} , the rule also covers assignments to those. It uses an auxiliary function $\overline{eval}_{\text{rexp}}$, which takes a reference expression (or \mathfrak{S}) re and an abstract

state $\bar{\sigma}$ and yields (a) a set NC of pairs, each consisting of an abstract reference to which re may evaluate in $\bar{\sigma}$ and the condition under which re may evaluate to this reference, and (b) a resulting abstract state, which is used to encode allocation, that is, when re contains new expressions (see Sec. 4.1).

$$(NC_{rhs}, (\bar{St}_{rhs}, \bar{con}_{rhs}, isSum_{rhs})) = \overline{eval}_{\text{rexp}}(rhs, \bar{\sigma}) \quad (1)$$

$$(NC_p, -) = \overline{eval}_{\text{rexp}}(p, (\bar{St}_{rhs}, \bar{con}_{rhs}, isSum_{rhs})) \quad (2)$$

$$strong \iff \exists \bar{n} \in \overline{Ref} \cdot (NC_p = \{(\bar{n}, -)\} \wedge \neg isSum_r(\bar{n})) \quad (3)$$

$$strong \Rightarrow (\bar{St} = \{(\bar{n}_1, f, \bar{n}_2) \in \bar{St}_{rhs} \mid (\bar{n}_1, -) \notin NC_p \vee f_r \neq f\}) \quad (4)$$

$$(\neg strong) \Rightarrow (\bar{St} = \bar{St}_{rhs}) \quad (5)$$

$$\bar{con}_{asg} = \left[(\bar{n}_p, f_r, \bar{n}_{rhs}) \mapsto (\text{Trg ToSrc}(\bar{s}_p) \sqcap \bar{s}_{rhs}) \mid \begin{array}{l} (\bar{n}_p, \bar{s}_p) \in NC_p \wedge \\ (\bar{n}_{rhs}, \bar{s}_{rhs}) \in NC_{rhs} \end{array} \right] \quad (6)$$

$$\bar{St}' = \bar{St} \cup \text{dom}(\bar{con}_{asg}) \quad (7)$$

$$\bar{con}' = \bar{con}_{rhs} \left[\bar{e} \mapsto \bar{s} \mid \begin{array}{l} \bar{e} \in \text{dom}(\bar{con}_{asg}) \wedge (\bar{e} \notin \bar{St} \Rightarrow \bar{s} = \bar{con}_{asg}(\bar{e})) \wedge \\ (\bar{e} \in \bar{St} \Rightarrow \bar{s} = \bar{con}_{asg}(\bar{e}) \sqcup \bar{con}_{rhs}(\bar{e})) \end{array} \right] \quad (8)$$

$$\langle p.f_r = rhs, \bar{\sigma} \rangle \rightarrow_{\Sigma} (\bar{St}', \bar{con}', isSum_{rhs})$$

A reference assignment first evaluates rhs to obtain the possible abstract references for the right-hand side expression together with the corresponding conditions, as well as a successor state (1). The receiver p is evaluated in this successor state. Since it is side-effect free (see Fig. 3), we discard the state resulting from its evaluation (2). The analysis performs a strong update iff there is only one abstract reference \bar{n} for the receiver, which is a definite node (3). In that case, the analysis removes all edges whose source is the receiver node and that are labeled with the assigned field f_r (4); otherwise, it performs a weak update, that is, retains all existing edges (5). To add the edges for all possible combinations of receivers and right-hand sides, we first create a map \bar{con}_{asg} that maps each of the new edges to the abstract condition that describes when the particular combination exists, that is, the greatest lower bound of the conditions for choosing a particular abstract reference for the receiver and a particular abstract reference for the right-hand side, respectively (6). The only twist in this step is how to handle edge-local identifiers. The receiver is denoted by **Trg** in conditions on edges pointing to the receiver, but by **Src** in the new edges. Function **Trg ToSrc** performs this conversion. Since the map \bar{con}_{asg} contains an entry for each new edge, we obtain the final abstract store by adding the domain of this map to the store constructed in step 4 or 5 (7). Finally, the abstract conditions are updated: For each new edge that is not present in the store before the reference assignment, we add the condition from \bar{con}_{asg} . For each edge that is already present (which may happen during a weak update), we join the condition from \bar{con}_{asg} and the existing condition (8).

Example. Fig. 9 *without* the bold edges and the highlighted constraints, shows the abstract heap after line 5 in Fig. 1. It is obtained from the abstract heap in Fig. 4 (the loop invariant) by (i) assuming the loop guard ($i > 0$) in all abstract conditions and (ii) applying the abstract semantics of the statement **Node p = new Node()** (line 5), which introduces the definite node \bar{n}_3 . (Its next

each way of evaluating rhs (1). Analogously to step 2 of reference assignment, we evaluate the receiver expression \mathbf{p} to obtain the possible receiver references, each with a condition under which \mathbf{p} may evaluate to this reference (2). We use the function $\text{update}_{\overline{\text{Cond}}}$ (see below) to reflect the value assignment in the value states of all edges in the abstract store (3). This function considers all possible combinations of receiver reference (obtained from $NC_{\mathbf{p}}$) and value state for a particular way of evaluating the right-hand side expression (from \overline{S}). For each of them, it propagates the value information that has to hold when this combination is chosen to the conditions of each edge in the abstract store and applies the assignment operation of the value domain. The condition of each edge in the abstract store is then defined to be the join of the conditions obtained for all ways of executing the value assignment.

Example. Fig. 9 *without* the highlighted constraints shows the abstract heap after line 6 in Fig. 1. The highlighted constraints are introduced by the abstract semantics of the statement $\mathbf{p.val} = \mathbf{i} - 1$ (line 7). There is only one way to evaluate the right-hand side expression. Therefore, $\text{eval}_{\text{vexp}}$ yields a singleton set (point (1) of the rule above). This set contains the condition $\mathbf{v} - \mathbf{i} \geq 0 \wedge \mathbf{i} > 0$, which holds in each concrete heap (otherwise there would be no value for result). Similarly, the receiver expression \mathbf{p} evaluates to a single node, \overline{n}_3 , under the same condition (2). This condition must be satisfied in order to be able to perform the assignment. Therefore, we conjoin it to each abstract condition in the store (which has no effect in this example), and then assign $\mathbf{i} - 1$ to $(\overline{n}_3, \text{val})$ since \mathbf{p} evaluates to \overline{n}_3 (3). Moreover, since \overline{n}_3 is the target of the edge from \mathbf{p} to \overline{n}_3 , we also add the constraint $(\text{Trg}, \text{val}) = \mathbf{i} - 1$ for the edge-local identifier to the condition on this edge, and analogously for (Src, val) on both out-edges of \overline{n}_3 (3).

Updating Conditions. Here, we formalize an update function used in above in the definition of the abstract semantics of value assignment.

We assume that the value analysis has the notion of summary identifiers and provides an operator $\text{assign} : \mathcal{P}(\overline{\text{ID}}) \times \text{vexp} \times \overline{\mathbf{V}} \rightarrow \overline{\mathbf{V}}$ that performs a strong assignment of a given expression to a given set of identifiers. We define a weak assignment $\text{assign}_w : \mathcal{P}(\overline{\text{ID}}) \times \text{vexp} \times \overline{\mathbf{V}} \rightarrow \overline{\mathbf{V}}$ as $\text{assign}_w(\overline{I}, \text{ve}, \overline{s}) = \text{assign}(\overline{I}, \text{ve}, \overline{s}) \sqcup \overline{s}$.

We define an update function $\text{update}_{\overline{\text{Cond}}} : (\text{vexp} \times \text{Field}_{\text{Val}} \times \mathcal{P}(\overline{\text{Ref}} \times \overline{\mathbf{V}}) \times \mathcal{P}(\overline{\mathbf{V}}) \times \overline{\text{Cond}}) \rightarrow \overline{\text{Cond}}$, which updates abstract conditions of all edges to reflect an assignment to a value field of a receiver with respect to the value states that correspond to the evaluations of the right hand side of the assignment. Formally:

$$\text{update}_{\overline{\text{Cond}}}(ve, f_v, NC, \overline{S}, \overline{\text{con}}) = \left[\overline{e} \mapsto \bigsqcup \left\{ \overline{s}_3 \left| \begin{array}{l} \exists \overline{s} \in \overline{S}, (\overline{n}, \overline{s}') \in NC. \\ \exists \overline{s}_1 = \text{propagate}_{\overline{\mathbf{V}}}(\overline{s}, \overline{s}', \overline{\text{con}}(\overline{e})). \\ \exists \overline{s}_2 = \text{asg}(\text{affected}_{\overline{\text{ID}}}(\overline{e}, f_v, \overline{n}), ve, \overline{s}_1). \\ \overline{s}_3 = \text{drop}(\text{vAE}, \overline{s}_2) \end{array} \right. \right\} \mid \overline{e} \in \text{dom}(\overline{\text{con}}) \right]$$

$$\text{where } \text{asg} = \begin{cases} \text{assign} & \text{if } \neg \text{isSum}(\overline{n}) \\ \text{assign}_w & \text{otherwise} \end{cases}.$$

For every edge \bar{e} , the function computes a new abstract condition as a join of conditions \bar{s}_3 computed as follows:

The set \bar{S} of value states corresponds to all possible evaluations of the value expression ve . NC is a set of abstract references representing a receiver together with an abstract condition under which this reference represents the receiver. A possible combination of \bar{s} in \bar{S} and (\bar{n}, \bar{s}') in NC is chosen. The value information from \bar{s} and \bar{s}' is propagated (via $propagate_{\bar{V}}$, defined below) into the abstract condition of \bar{e} , obtaining the value state \bar{s}_1 . \bar{s}_1 contains all value field access identifiers from ve . Since several different identifiers in \bar{ID} may represent the same memory location, all the identifiers that may represent the value field f_v of the receiver \bar{n} must be updated. We define a function $affected_{\bar{ID}} : ((\bar{Ref} \times \text{Field}_{\bar{Ref}} \times \text{Ref}) \times \text{Field}_{\text{Val}} \times \text{Ref}) \rightarrow \mathcal{P}(\bar{ID})$ which yields the set of identifiers for a given edge that may represent the same memory location as the identifier (\bar{n}, f_v) . Formally:

$$affected_{\bar{ID}}((\bar{n}_1, f_r, \bar{n}_2), f_v, \bar{n}) = \{(\text{Src}, f_v) \mid \bar{n}_1 = \bar{n}\} \cup \{(\text{Trg}, f_v) \mid \bar{n}_2 = \bar{n}\} \cup \{(\bar{n}, f_v)\}$$

Then, we apply a strong (if \bar{n} is definite) or a weak (if \bar{n} is summary) assignment of ve to all identifiers that may represent the same memory location as (\bar{n}, f_v) , obtaining the value state \bar{s}_2 . However, \bar{s}_2 still contains value field access identifiers of ve which cannot appear in abstract conditions of edges. We obtain \bar{s}_3 by removing all value field access identifiers from \bar{s}_2 (via $drop$).

Function $propagate_{\bar{V}} : \bar{V} \times \bar{V} \times \bar{V} \rightarrow \bar{V}$ propagates the value information from a given value state \bar{s} corresponding to an evaluation of a value expression and a given abstract condition \bar{s}' for a receiver to a given abstract condition of an edge \bar{s}_e . Formally:

$$propagate_{\bar{V}}(\bar{s}, \bar{s}', \bar{s}_e) = meet_{\text{VAE}}(add(\bar{EId}, meet_{\text{VAE}}(\bar{s}, \bar{s}')), \bar{s}_e)$$

The above function first gathers the value information from \bar{s} and \bar{s}' into an intermediate state by taking the greatest lower bound of \bar{s} and \bar{s}' and preserving the value field access identifiers of the both states. Then, in order to preserve the edge local identifiers of \bar{s}_e , the edge-local identifiers are introduced into the intermediate state. Finally, we propagate the value information from the intermediate state to \bar{s}_e by taking the greatest lower bound of the intermediate state and \bar{s}_e while preserving the value field access identifiers of the intermediate state.

6 Experimental Results

We implemented our analysis in the static analyzer Sample and applied it to Scala implementations of typical list and tree operations (some of which we took from the literature [4,6,13]), operations on nested recursive data structures (such as lists of lists), and a simple aggregate structure [11]. We performed the experiments on an Intel Core i7-Q820 CPU (1.73GHz, 8GB) running the 64-bit version of Ubuntu 14.04. We instantiated our analysis with the Octagon [22]

Data Structure	Operations	Octa.	Poly.	Data Structure	Operations	Octa.	Poly.	
SortedSLL	constructor	1.24	1.82	BST	constructor	1.96	2.43	
	insertKey				insertKey			
	deleteKey			findKey	NodeCachingSLL	constructor	0.87	1.04
	findKey			deepCopy		add		
SortedDLL	constructor	1.91	2.83	PersonAndAccount	remove	0.38	0.43	
	insertKey				findKey			
	deleteKey				deepCopy			findKey
	deepCopy				withdraw			
					deposit			
					changeInterest			

Table 1: Analysis times (in seconds) of classes implementing different data structures when instantiating the analysis with the Octagon and Polyhedra value domains. For each class, we inferred an object invariant by computing a fixed point over all its methods.

and Polyhedra [8] value domains implemented in Apron [16]. We used the default widening parameter, that is, \mathcal{W} is the set of local reference variables. There were no manual annotations for any of the benchmarks.

Inference of Object Invariants. Tab. 1 reports the analysis times (the average of 10 runs) for implementations of five different data structures. We instantiated Logozzo’s framework [18] with our analysis to infer object invariants for each data structure by computing a fixed point over all its operations.

SortedSLL implements a sorted singly-linked list (SLL). The inferred object invariant expresses that the values stored in the list nodes are non-decreasing.

SortedDLL implements a sorted doubly-linked list (DLL). Our analysis infers sortedness in both directions, that is, via the `next` and `prev` fields. However, the analysis cannot infer the structural invariant of doubly-linked lists $n.\text{next} \neq \text{null} \Rightarrow n = n.\text{next}.\text{prev}$ because it has no way of relating the concrete references of the two edges $\bar{n} \xrightarrow{\text{next}} \bar{n}$ and $\bar{n} \xrightarrow{\text{prev}} \bar{n}$ for the summary node \bar{n} .

BST implements a binary search tree. The analysis infers both the value and the shape information of a BST data structure. Our implementation stores the infimum and supremum of all keys of a sub-tree in its root. This information allows our analysis to distinguish the left and right sub-tree of a node and, thus, to infer that the shape is not a DAG. We omitted method `deleteKey` because our analysis is not able to infer that replacing the deleted key with the next smallest key preserves sortedness; it does, however, infer that the tree shape is preserved.

NodeCachingSLL implements an acyclic SLL that maintains a cache of node objects to reduce object creation and garbage collection. The inferred object invariant expresses that the list and the cache are disjoint and that the size of the cache is between 0 and `maximumCacheSize`. Moreover, we inferred that the `addKey` method creates a new object only if the cache is empty. Every node of the list stores the length of the list rooted at the node. This information lets our analysis infer that the list and its cache are acyclic, which is needed to infer disjointness of the list and the cache. The latter step required materialization, that is, splitting a definite node off a summary node, which is supported by our implementation, but not explained in this paper.

Operation	Octa.	Poly.
insertionSort	0.43 - SLL	0.48 - SLL
	0.72 - DLL	0.85 - DLL
partitionWithKey	0.32 - SLL	0.34 - SLL
	0.48 - DLL	0.55 - DLL
createListOfZerosAndSum	0.22 - SLL	0.23 - SLL
	0.39 - DLL	0.43 - DLL
increasingList	0.28 - SLL	0.31 - SLL
	0.41 - DLL	0.50 - DLL
sortListOfListsOfValues	1.45	1.88
listToBST	1.03	1.21

Table 2: Analysis times (in seconds) for single operations on different data structures when instantiating the analysis with the Octagon and Polyhedra value domains. The first 4 operations work on singly-linked lists (SLL) and doubly-linked lists (DLL). The fifth operation works on lists of singly-linked lists that store values. The last operation transforms an SLL to a binary search tree.

Besides the object invariants for these four classes, our analysis infers that the result of method `findKey` is either null or has the value of the given key. This postcondition is inferred even if the result is represented by a summary node.

`PersonAndAccount` implements an aggregate data structure similar to the example from a paper on the verification of object invariants [11]. The analysis infers combined shape and value invariants, for instance, that `Account` and `Person` objects reference each other, the sum of the account balance and person’s salary is positive, and the interest rate of the account is always non-negative.

Inference of Method Postconditions. Tab. 2 reports the analysis times of individual operations on different data structures. The initial abstract states and the abstract heaps that represent the arguments to the operations contain only information that is provided by the static types; no annotations were used. The first four operations manipulate singly and doubly-linked lists. `insertionSort` takes an unsorted list of values and sorts it. The analysis infers that the result is a sorted list. `partitionWithKey` takes a list of values and a key, and creates two new lists such that the keys in one are less than or equal to the given key, and the keys in the other are greater. The analysis infers this value property and that the resulting lists are disjoint. `createListOfZerosAndSum` creates a list of zeros and subsequently traverses the list and sums up the values. The analysis infers that the result is a list of zeros, and the sum of the values is zero. `increasingList` is the method from Fig. 1, with an analogous implementation for DLLs. The analysis infers the heap in Fig. 2 (and an analogous heap for DLLs).

The last two operations of Tab. 2 demonstrate that the analysis is able to infer non-trivial shape and value properties for programs manipulating nested recursive data structures or a combination of different data structures. `sortListOfListsOfValues` takes a singly-linked list of SLLs that store values, and sorts each of the lists. The analysis infers that the result is a list of sorted SLLs. `listToBST` takes an SLL of values and creates a binary search tree out of it, without using the methods of the `BST` class discussed above. The analysis infers that the result is a binary search tree.

Discussion. The analysis times in Tab. 1 and Tab. 2 demonstrate the efficiency of our analysis. For all our benchmark classes, the fixed point over all their methods was computed within 3 seconds when using the Polyhedra domain. When instantiated with a more efficient but less precise value domain, the efficiency of the analysis increases, as illustrated by the usage of the Octagon domain.

Our experiments demonstrate that our analysis can infer invariants that combine shape and value information in interesting ways, for instance, sortedness of lists and trees, or invariants that relate the states of different objects in an aggregate structure. Our analysis leverages data stored in value fields, such as the infimum and supremum in the `BST` class discussed above, to obtain more precise shape information. As future work, we plan to rely less on such fields by tracking additional abstract conditions (such as injectivity of references) on edges and by generalizing edge-local identifiers to reference fields.

7 Related Work

Dillig et al. [9,10] present a precise content analysis for arrays and containers, in which heap edges are qualified by logical constraints over indexes into a container. This idea inspired our approach of tracking disjunctive information via the value states associated with edges in the heap. Our analysis uses generic value domains instead of logical constraints and can therefore be instantiated with different levels of precision and efficiency. Moreover, it uses edge-local identifiers instead of indexes, which allows us to express constraints on arbitrary nodes (especially summary nodes) in the heap, not only on indexed structures such as arrays and containers. Whereas Dillig et al. concentrate on clients of arrays and containers, our analysis targets arbitrary heap-manipulating programs including implementations of containers.

Similarly to our work, Bouajjani et al. [1,2] introduce a static analysis that automatically infers combined shape and numerical invariants and is parametric in the underlying value domain. The main difference is that their technique is specific to programs that manipulate singly-linked lists of values. For such data structures Bouajjani et al.’s approach is more powerful since it can relate an arbitrary number of successive positions in a list. In contrast, the aim of our analysis is to be applicable to general heap-manipulating programs.

Sagiv et al. [23] introduce a shape analysis in which invariants are expressed in 3-valued first order logic with transitive closure (FOLTC). These invariants may combine shape and value constraints. The analysis requires user-supplied predicates, whereas our analysis does not need manual annotations; it represents a state by a set of logical structures, whereas our analysis maintains a single abstract heap, reducing the number of nodes and edges, and therefore the complexity of the overall analysis. The merging of nodes in our widening operator can be viewed as a special case of canonical abstractions.

McCloskey et al. [21] propose a framework for combining shape and numerical domains (encoded as predicates in FOLTC) in a generic way. However, users have to supply shared predicates via which the domains communicate and which

usually resemble the properties one wants to prove. In contrast, our analysis can be parameterized by arbitrary value domains without any manual overhead.

Ferrara et al. [12,13] and Fu [14] combine different heap and value analyses. Whereas their work represents a state as a heap abstraction and a single value state, our analysis attaches a value state to each edge in the heap abstraction, allowing for a precise tracking of disjunctive information. Moreover, in the value states of Ferrara et al.’s and Fu’s work, different heap identifiers represent disjoint portions of the heap. This is not the case for our edge-local identifiers, which refer to memory locations already represented by abstract identifiers and which enable a precise treatment of summary nodes.

Chang et al. [6] introduce a shape analysis based on user-supplied invariants that describe data structures such as lists and trees. These invariants are used to abstract over a potentially unbounded number of concrete references. Chang and Rival [4,5] extend this work and present a framework for combining shape and numeric abstractions into a single domain. Their approach enables the precise and modular analysis of heap and numeric invariants, but relies on user-supplied properties, whereas our analysis does not require manual annotations.

Marron et al. [20] introduce heap abstractions that are similar to the graphs representing abstract heaps in our work. In fact, the formalization of the concretization function in Sec. 3.2 is inspired by their work. However, there are important technical differences. In particular, Marron et al.’s analysis maintains a normal form, which makes their lattice finite, but loses information when merging two heap graphs. In contrast, we deal with an infinite lattice, but preserve some of this information. Moreover, Marron et al.’s heap graphs track specific aliasing predicates (such as injectivity of fields or tree shapes), but no value information. Finally, the purpose of their work is to provide a high-level abstraction of *concrete* runtime heaps, whereas we propose an abstract domain and an abstract semantics for a static code analysis.

8 Conclusion

In this paper, we have presented a static analysis that infers complex invariants combining shape and value information. The analysis is parametric in the underlying value domain, allowing for different trade-offs between precision and efficiency. A key innovation of our analysis is the introduction of edge-local identifiers to track value information about the source and target of a pointer, which allows it to infer inductive invariants such as sortedness of a linked list. The analysis has been implemented in the static analyzer *Sample*. Our experiments demonstrate its effectiveness.

Acknowledgments. We are grateful to Uri Juhasz and Alexander Summers for numerous discussions, to John Boyland for helpful comments on a draft of this paper, and to Severin Heiniger for his contributions to the implementation.

References

1. A. Bouajjani, C. Dragoi, C. Enea, and M. Sighireanu. On inter-procedural analysis of programs with lists and data. In *PLDI*. ACM, 2011.
2. A. Bouajjani, C. Dragoi, C. Enea, and M. Sighireanu. Abstract domains for automated reasoning about list-manipulating programs with infinite data. In *VMCAI*. Springer, 2012.
3. H. Bunke, X. Jiang, and A. Kandel. On the minimum common supergraph of two graphs. *Computing*, 65(1):13–25, Aug. 2000.
4. B.-Y. E. Chang and X. Rival. Relational inductive shape analysis. In *POPL*. ACM, 2008.
5. B.-Y. E. Chang and X. Rival. Modular construction of shape-numeric analyzers. In *David A. Schmidt’s 60th Birthday Festschrift*, EPTCS, 2013.
6. B.-Y. E. Chang, X. Rival, and G. C. Necula. Shape analysis with structural invariant checkers. In *SAS*. Springer, 2007.
7. P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *POPL*. ACM, 1977.
8. P. Cousot and N. Halbwachs. Automatic discovery of linear restraints among variables of a program. In *POPL*. ACM, 1978.
9. I. Dillig, T. Dillig, and A. Aiken. Fluid updates: Beyond strong vs. weak updates. In *ESOP*. Springer, 2010.
10. I. Dillig, T. Dillig, and A. Aiken. Precise reasoning for programs using containers. In *POPL*. ACM, 2011.
11. S. Drossopoulou, A. Francalanza, P. Müller, and A. J. Summers. A unified framework for verification techniques for object invariants. In *ECOOP*. Springer, 2008.
12. P. Ferrara. Generic combination of heap and value analyses in abstract interpretation. In *VMCAI*. Springer, 2014.
13. P. Ferrara, R. Fuchs, and U. Juhász. TVAL+ : TVLA and value analyses together. In *SEFM*. Springer, 2012.
14. Z. Fu. Modularly combining numeric abstract domains with points-to analysis, and a scalable static numeric analyzer for Java. In *VMCAI*. Springer, 2014.
15. D. Gopan, F. DiMaio, N. Dor, T. W. Reps, and S. Sagiv. Numeric domains with summarized dimensions. In *TACAS*. Springer, 2004.
16. B. Jeannet and A. Miné. Apron: A library of numerical abstract domains for static analysis. In *CAV*. Springer, 2009.
17. T. Lev-Ami and M. Sagiv. TVLA: A system for implementing static analyses. In *SAS*. Springer, 2000.
18. F. Logozzo. Automatic inference of class invariants. In *VMCAI*. Springer, 2004.
19. R. Manevich, M. Sagiv, G. Ramalingam, and J. Field. Partially disjunctive heap abstraction. In *SAS*. Springer, 2004.
20. M. Marron, C. Sánchez, Z. Su, and M. Fähndrich. Abstracting runtime heaps for program understanding. *IEEE Trans. Software Eng.*, 39(6):774–786, 2013.
21. B. McCloskey, T. W. Reps, and S. Sagiv. Statically inferring complex heap, array, and numeric invariants. In *SAS*. Springer, 2010.
22. A. Miné. The octagon abstract domain. *Higher Order Symbol. Comput.*, 2006.
23. M. Sagiv, T. Reps, and R. Wilhelm. Parametric shape analysis via 3-valued logic. In *POPL*. ACM, 1999.

A Graph Theory

Definition 1 (Graph). Let L be a finite set of labels of nodes and edges. A labeled **graph** is a triple $g = (V, E, \eta, \mu)$, where

- V is a set of vertices
- $E \subseteq V \times L \times V$
- $\eta : V \rightarrow L$ is a function assigning labels to vertices

Definition 2 (Subgraph - \subseteq). Let $g = (V, E, \eta)$ be a graph. We say that $g' = (V', E', \eta')$ is a **subgraph** of g , $g' \subseteq g$, iff

- $V' \subseteq V$
- $E' \subseteq E$
- $\eta'(v') = \eta(v')$ for all $v' \in V'$

Definition 3 (Graph Isomorphism). Let g and g' be graphs. A **graph isomorphism** between g and g' is a bijective mapping $i : V \rightarrow V'$ such that

- $\eta(v) = \eta'(i(v))$ for all $v \in V$
- for any edge $e = (u, l, v) \in E$ there exists an edge $e' = (i(u), l, i(v)) \in E'$, and for any edge $e' = (u', l, v') \in E'$ there exists an edge $e = (i^{-1}(u'), l, i^{-1}(v')) \in E$.

Definition 4 (Subgraph Isomorphism). If there exists a graph $g' \subseteq g$ and g' is isomorphic to g_1 , then we say that g_1 is subgraph isomorphic to g .

Definition 5 (Maximum Common Subgraph - mcs). Let $g_1 = (V_1, E_1, \eta_1)$ and $g_2 = (V_2, E_2, \eta_2)$ be graphs. We say that $g = (V, E, \eta)$ is a **common subgraph** of g_1 and g_2 iff there exist $g'_1 \subseteq g_1$ and $g'_2 \subseteq g_2$ such that g is graph isomorphic to g'_1 and g'_2 . We call g a **maximum common subgraph** of g_1 and g_2 , $mcs(g_1, g_2)$, if there exists no other common subgraph of g_1 and g_2 that has, firstly, more vertices and, secondly, more edges than g .