



**HAL**  
open science

# Intraprocedural Pointer Analysis for Container-Centric Applications

Albert Cohen, Peng Wu, David Padua

► **To cite this version:**

Albert Cohen, Peng Wu, David Padua. Intraprocedural Pointer Analysis for Container-Centric Applications. RR-4289, INRIA. 2001. inria-00072298

**HAL Id: inria-00072298**

**<https://inria.hal.science/inria-00072298>**

Submitted on 23 May 2006

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

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

# *Intraprocedural Pointer Analysis for Container-Centric Applications*

Albert Cohen — Peng Wu — David Padua

**N° 4289**

Octobre 2001

THÈME 4



*R*  
*apport*  
*de recherche*



# Intraprocedural Pointer Analysis for Container-Centric Applications

Albert Cohen<sup>\*</sup>, Peng Wu<sup>†</sup>, David Padua<sup>†</sup>

Thème 4 — Simulation et optimisation  
de systèmes complexes  
Projet A3

Rapport de recherche n° 4289 — Octobre 2001 — 32 pages

**Abstract:** As programmers look forward to designing high performance applications with object-oriented models, compilers must support higher-level analyses and optimizations. Pointer analysis for container-centric applications is one of these: it exploits abstract semantics of container structures (e.g., lists, trees, associative maps) provided by standard libraries and toolkits. Extending shape analysis work by Sagiv, Reps and Wilhelm, we capture aliasing properties through dedicated points-to graphs. Formalization in abstract interpretation allowed us to prove the abstraction's and transfer functions' safety. We ran the analysis on small examples. It achieved precise memory disambiguations useful to parallelization and optimization.

**Key-words:** points-to graph, abstract container type, abstract interpretation, alias and shape analysis

<sup>\*</sup> A3 Project, INRIA Rocquencourt

<sup>†</sup> DCS, University of Illinois at Urbana-Champaign

# Analyse de pointeurs intraprocédurale pour les applications manipulant des conteneurs

**Résumé :** Les programmeurs s'intéressent de plus en plus aux modèles orientés-objets pour la conception d'applications à hautes performances ; les compilateurs doivent ainsi réaliser des analyses et des optimisations de plus haut niveau. Une de ces analyses s'intéresse aux pointeurs dans les applications manipulant des conteneurs : elle exploite la sémantique abstraite des structures de conteneurs (listes, arbres, tables associatives...) proposées par des bibliothèques et des boîtes à outils standard. En étendant l'analyse de forme (*shape analysis*) de Sagiv, Reps et Wilhelm, nous décrivons les relations d'alias à l'aide de graphes *points-to* spécifiques. Une preuve de correction des abstractions et des fonctions de transfert a été réalisée dans le formalisme de l'interprétation abstraite. L'analyse a été expérimentée sur de courts exemples. Elle permet de séparer précisément les références mémoire, mettant ainsi en évidence une parallélisation et des optimisations intéressantes.

**Mots-clés :** graphe *points-to*, type abstrait de conteneur, interprétation abstraite, analyse d'alias et de forme

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Overview of the Method</b>	<b>4</b>
<b>3</b>	<b>The Direction Abstraction</b>	<b>5</b>
3.1	Connections and Directions . . . . .	5
3.2	Abstract Interpretation for Directions . . . . .	7
<b>4</b>	<b>Points-to Graphs</b>	<b>8</b>
4.1	The Node Naming Scheme . . . . .	8
4.2	Definition of Points-to Graphs . . . . .	9
<b>5</b>	<b>Points-to Analysis</b>	<b>10</b>
5.1	An Intuitive Flavor of the Analysis . . . . .	10
5.2	Strong Nullification of Variables . . . . .	12
5.3	Node Materialization . . . . .	13
5.4	Transfer Functions . . . . .	13
5.4.1	The Classical Part . . . . .	13
5.4.2	Iterator Operations . . . . .	14
5.4.3	Container Operations . . . . .	15
5.5	Merging Edges . . . . .	16
5.6	Correctness Issues . . . . .	17
<b>6</b>	<b>Illustrative Examples</b>	<b>18</b>
6.1	List Copying and Reversal . . . . .	18
6.2	Swapping Elements . . . . .	18
6.3	In-Place Container Rotation or Reversal . . . . .	20
6.4	Nested Traversals . . . . .	21
<b>7</b>	<b>Aliasing and Combness</b>	<b>22</b>
<b>8</b>	<b>Putting the Analysis to Work</b>	<b>23</b>
8.1	Cost Improvements . . . . .	23
8.2	Precision Improvements . . . . .	24
<b>9</b>	<b>Related Work</b>	<b>24</b>
<b>10</b>	<b>Conclusion and Future Work</b>	<b>25</b>
<b>A</b>	<b>Working With Directions</b>	<b>26</b>
<b>B</b>	<b>Abstract Interpretation Proofs</b>	<b>28</b>
B.1	Combness and Comparison . . . . .	28
B.2	Simplifications . . . . .	28
B.3	Distributivity and Partial Distributivity . . . . .	29
B.4	Lattice Completeness . . . . .	30
B.5	Galois Connection . . . . .	30
B.6	Local Safety . . . . .	31
B.7	Merging Edges . . . . .	32

# 1 Introduction

Object-oriented design plays an increasing role in performance-critical codes. The shift toward higher level design models enables software reuse, generic programming, and aids the construction of large software systems. When dealing with general-purpose programs, arrays should share their preeminence with more general *container* components, such as lists, hash-tables, sets, and so on. These general-purpose containers are often provided by standard libraries, such as C++ Standard Template Library (STL) and Java Standard Development Kit (SDK).

This work aims at the automatic retrieval of pointer properties in the context of container classes. Besides classical points-to relations, we are especially interested in the *connection* between containers and their elements. Several intuitions were borrowed from [23]: pointer analysis should not limit itself to *pointer-chasing*, i.e., recursive dereferences of the form `p = p.next` (in Java syntax). Indeed, object-oriented design tends to hide such technicalities behind *containers* and *iterators* that traverse them. Although implementations of container classes could be very complex, especially with hand optimizations and generic codings, they often exhibit clean semantics through their programming interfaces. Thus, for a better analyzability of object-oriented codes, the key is to exploit semantics of critical classes in compilers.

The paper presents a points-to analysis for container traversals. The presentation uses Java syntax and semantics (Java references are indeed “well-behaved” pointers), but the general framework is not limited to Java.<sup>1</sup> Studying real-world kernels, we show that our technique achieves a very high precision as soon as container traversals expose enough regularity. For some of the kernels, the same precision cannot be achieved by classical techniques applied to an equivalent pointer-chasing implementation. The result of the analysis can be interpreted as shape [11, 19], alias [3, 8] or connectivity [10]. These properties can further improve analyses such as dependence tests and various optimizations.

We also address cost issues of such a high-precision analysis. In fact, when efficiency matters, one may still be very precise about container operations while applying inexpensive rules on other pointer assignments. This is a clear benefit of using standard containers as opposed to hand-written pointer-chasing implementations: heuristics can be made conscious of the level of precision required, without additional “hints” from the programmer or the “blind” degradation of *k*-limiting.

Section 2 presents the motivation for this work and introduces the main concepts. Section 4 describes the points-to graph and how to abstract the connection between a container and its elements. The points-to algorithm is given in Section 5. More general alias and shape analysis applications are presented in Section 7. Section 6 applies our technique to illustrative examples. Section 8 discusses possible improvements. We then compare our approach with others in Section 9. Section 10 concludes and sketches future work. All technical proofs have been moved to Section B in the appendix.

## 2 Overview of the Method

The analysis targets containers that can be traversed as linear lists, e.g., lists and vectors. We only studied a small set of basic container operations: `first()` and `last()` generates a new iterator attached to the container; `advance()` and `retreat()` move an iterator; `get()`, `put()`, `insert()` and `delete()` update/accesses a container through an iterator. The syntax of these operations is borrowed from the Java Generic Library (JGL).<sup>2</sup> The semantics is slightly simplified, e.g., no implicit iterator move in container accesses.

Let’s consider the example in Figure 1 and suppose that `list` refers to a *non-circular* list. We would like to infer whether `a` and `b` are aliased at point 6. This would depend on two conditions: whether iterator `i` indicates distinct positions at points 2 and 4; and, whether distinct objects are attached to distinct positions inside `list`. The first condition can be decided by tracking “moves” of an iterator. The second condition is related to the “shape” of a container. A container is a *comb* when *distinct elements are attached at distinct positions*. *Combness* is the container counterpart of classical shape attributes, e.g., *listness* or *treeness*. Figure 1 gives both examples of `list` pointing to a *comb* or *not comb* container: `a` and `b` are unaliased in the first graph, and aliased in the second one. Moreover, after executing `i.put(o)`, one would like to replace

---

<sup>1</sup>Extension to C/C++ should be possible using well-known techniques [16, 21, 2].

<sup>2</sup>A project from ObjectSpace, see <http://www.objectspace.com>.

```

1 Foo o = new Object ();
2 Iterator i = list.first ();
3 Foo a = i.get ();
4 i.advance ();
5 Foo b = i.get ();
6 i.put (o);

```

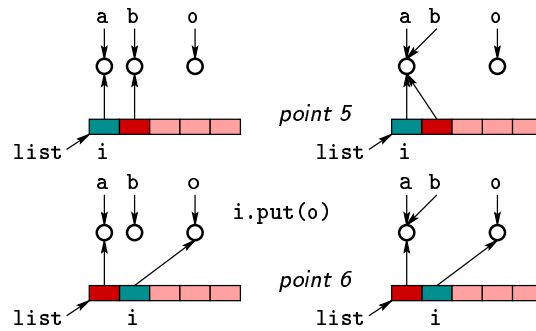


Figure 1: Motivating example

the element referenced by `i` with the object referenced by `o`. This implies killing the connection between `list` and the object referenced by `b`.

In this paper, combness is used as an intermediate property for capturing points-to, aliasing and connectivity relations. Nevertheless, it might also be applied directly to parallelization (combined with dependence analysis) [23] and data structure layout optimizations, such as object inlining [9] (combness is similar to so called, one-to-one fields) and “flattening” of Java arrays [22].

Overall, the analysis is based on a graph representation that abstract the points-to relation between pointers. Our technique is strongly influenced by Sagiv, Reps and Wilhelm’s shape analysis for destructive updating [18]. In addition, container-specific information is encoded into the graph, such as iterator attachment, combness and container positions. Focusing on precision, the analysis is flow-sensitive and handles destructive updates of the form `v.f = null` (for general pointers) or `i.delete()` (for containers). We also introduce a new scheme to abstract element retrieval through `get()`. It is a generalization of Sagiv, Reps and Wilhelm’s *materialization* [18].

### 3 The Direction Abstraction

The points-to analysis needs to model precisely the connection between containers and the objects they store. This section introduces the *direction* abstraction to capture properties of such connections.

#### 3.1 Connections and Directions

Our points-to graphs uses special container edges to model the link between containers and the objects attached. We call the concrete state abstracted by such container edge a *connection*. In particular, objects abstracted by the source/target nodes of the edge are called the *source containers/target objects* of the connection. A connection has two important properties, the *domain* and the *injectivity*,

- The *domain* of a connection describes at which positions the target objects of the connection are attached at its source containers.
- The *injectivity* of a connection tells whether each target object of the connection is attached at a single position of the source containers.

We use *directions* to describe the domain and injectivity of a connection. In the points-to graph, this is reflected as labeling container edges by directions. First, let us define *primitive directions*, examples are given in Figure 2.

- $\diamond_i$  describes any injective connection whose domain is the *single* position of (indicated by) `i`;
- $\triangleleft_i$  describes any injective connection whose domain *precedes* the position of `i`;
- $\triangleright_i$  describes any injective connection whose domain *follows* the position of `i`;
- $\emptyset$  describes any connection whose domain is empty;



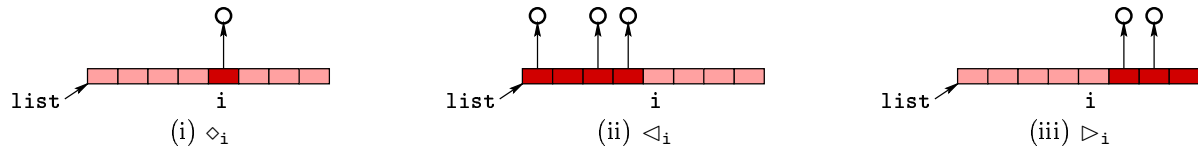


Figure 2: Connections described by primitive directions

\* describes any injective connection.

To abstract more properties of connections, we define *direction expressions*. A *direction expression* can be a primitive direction, or it can be built from the four operators specified below. In the following, we say a connection *satisfies* a direction if it can be described by the direction. Examples of direction expressions are given in Figure 3.

- $e \cdot e'$  describes any connection that satisfies both  $e$  and  $e'$ . Operator  $\cdot$  is called the *meet* operator. When clear of the context,  $\cdot$  is left out.
- $e \# e'$  describes any connection that can be split into two connections, such that one satisfies  $e$  and the other satisfies  $e'$ . Operator  $\#$  is called the binary *mix* operator.
- $e \parallel e'$  describes any connection that can be split into two connections with *disjoint* target objects such that one satisfies  $e$  and the other satisfies  $e'$ . Operator  $\parallel$  is called the *join* operator.
- $e^\#$  describes any connection whose domain satisfies  $e$ . Operator  $\#$  is called the unary *mix* operator.

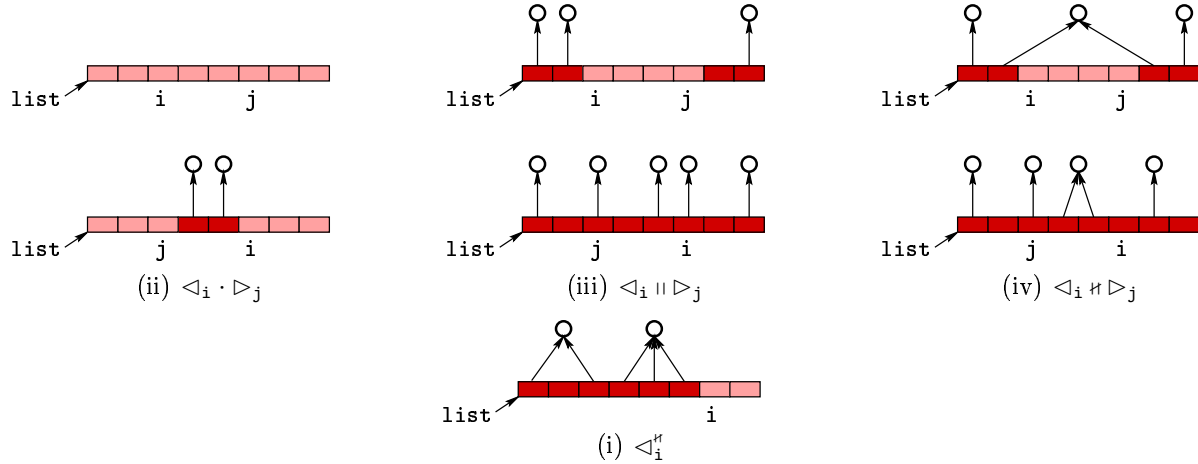


Figure 3: Examples of Direction Expressions

Finally, let  $\text{Expr}$  be the set of all direction expressions. We define an equivalence relation  $\equiv$  on  $\text{Expr}$ :  $e \equiv e'$  if and only if the same set of connections satisfy  $e$  and  $e'$ . An equivalence class for  $\equiv$  is called a *direction*. The set of all directions,  $\text{Dir}$ , is thus the *quotient* of  $\text{Expr}$  by  $\equiv$ , i.e.,  $\text{Expr}/\equiv$ .

According to the above definitions, injectivity of connections can be inductively deduced from direction expressions. More precisely, a direction is a *comb* if it describes *injective* connections only. Combness of a direction  $d$ , written  $\text{Comb}(d)$ , can be computed inductively from Table 1 (proof in Section B.1).

$\text{Comb}(d \cdot d') = \text{Comb}(d) \vee \text{Comb}(d')$ $\neg \text{Comb}(d \# d')$ , (generalcase) $\neg \text{Comb}(d^\#)$ , (generalcase)	$\text{Comb}(d \parallel d') = \text{Comb}(d) \wedge \text{Comb}(d')$
------------------------------------------------------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------

Table 1: Combness of direction expressions

Combsness of directions improve the precision of the points-to analysis: from the injectivity of connections, one may deduce the *shape* of a container. On the other hand, we also need to tell whether an object is attached at a particular container position. Obtaining the latter information involves comparing the domains of two connections. Therefore, we define the following relations on directions:

- *Disjointness* tells whether two directions describe connections with disjoint domains;
- *Coincidence* tells whether two directions describe connections whose domains are *singletons* holding the *same* position.

Table 2 defines predicates *Coincide* and *Disjoint* which test these relations. Notice *Disjoint* is not the complement of *Coincide*. For instance, given  $d = (\diamond_i \triangleleft_j) \uplus \diamond_k$  and  $d' = \diamond_i \uplus (\diamond_k \triangleright_j)$ , neither *Coincide*( $d, d'$ ) nor *Disjoint*( $d, d'$ ) hold.

Coincide( $\diamond_i, \diamond_i$ )	Disjoint( $\triangleleft_i, \diamond_i$ ), Disjoint( $\diamond_i, \triangleright_i$ ), Disjoint( $\triangleleft_i, \triangleright_i$ )
Coincide( $d, d_1 \cdot d_2$ ) = Coincide( $d, d_1$ ) $\vee$ Coincide( $d, d_2$ )	Disjoint( $d, d_1 \cdot d_2$ ) = Disjoint( $d, d_1$ ) $\vee$ Disjoint( $d, d_2$ )
Coincide( $d, d_1 \uplus d_2$ ) = Coincide( $d, d_1$ ) $\wedge$ Coincide( $d, d_2$ )	Disjoint( $d, d_1 \uplus d_2$ ) = Disjoint( $d, d_1$ ) $\wedge$ Disjoint( $d, d_2$ )
Coincide( $d, d_1 \uplus d_2$ ) = Coincide( $d, d_1$ ) $\wedge$ Coincide( $d, d_2$ )	Disjoint( $d, d_1 \uplus d_2$ ) = Disjoint( $d, d_1$ ) $\wedge$ Disjoint( $d, d_2$ )
Coincide( $d, d_2^{\#}$ ) = Coincide( $d, d_2$ )	Disjoint( $d, d_2^{\#}$ ) = Disjoint( $d, d_2$ )

Table 2: Coincidence and disjointness of directions

### 3.2 Abstract Interpretation for Directions

We use *abstract interpretation* [5] to enable systematic proofs about directions. *lter* is a finite set of iterator variables. Formally, a connection is a pair of mappings  $(\mu, \nu)$ , where  $\mu : \mathbb{N} \rightarrow \mathbb{N}$  maps container positions to the objects attached at the positions, and  $\nu : \text{lter} \rightarrow \mathbb{N}$  maps iterator variables to the positions they indicate.<sup>3</sup> Let *Conn* be the set of all connections.

- First, we define function  $\gamma : \text{Expr} \rightarrow 2^{\text{Conn}}$  for direction expressions, as given in Table 3. Directions are equivalence classes of *Expr* for the following relation:

$$\forall e, e' \in \text{Expr} : e \equiv e' \stackrel{\text{def}}{\iff} \gamma(e) = \gamma(e').$$

From this relation, we can define the *concretization* function  $\gamma : \text{Dir} \rightarrow 2^{\text{Conn}}$ . Then, the link with the intuitive definition of directions is the following: a connection,  $c$ , satisfies a direction,  $d$ , if and only if  $c \subseteq \gamma(d)$ .

- We then define a *partial ordering*  $\sqsubseteq$  over *Dir*:

$$d \sqsubseteq d' \stackrel{\text{def}}{\iff} \gamma(d) \subseteq \gamma(d').$$

One may easily check the following orderings on directions:

$$d'' \sqsubseteq d \cdot d' \implies d'' \sqsubseteq d \wedge d'' \sqsubseteq d' \quad d \uplus d' \sqsubseteq d'' \implies d \sqsubseteq d'' \wedge d' \sqsubseteq d''.$$

As a result, *Dir* can be structured as a complete lattice ordered by  $\sqsubseteq$ . The *join* and *meet* operators of the lattice are  $\uplus$  and  $\cdot$ , respectively. We prove the completeness of the lattice in Section B.4.

- Finally, the *abstraction* function  $\alpha : 2^{\text{Conn}} \rightarrow \text{Dir}$  is built from  $\gamma$ . Given  $C \subseteq \text{Conn}$ ,  $\alpha$  is defined as,

$$\alpha(C) = \min_{\sqsubseteq} \{d \in \text{Dir} : C \subseteq \gamma(d)\}.$$

The minimum exists because the lattice of *Dir* is complete. Indeed, an equivalent definition for  $\alpha(C)$  is the (infinite) meet of all directions that  $C$  satisfies.

$\gamma(*)$	$= \{(\mu, \nu) \mid \text{Injective}(\mu)\}$
$\gamma(\emptyset)$	$= \{(\mu, \nu) \mid \text{Domain}(\mu) = \emptyset\}$
$\gamma(\diamond_i)$	$= \{(\mu, \nu) \mid \text{Domain}(\mu) = \{\nu(i)\}\}$
$\gamma(\triangleleft_i)$	$= \{(\mu, \nu) \mid (\forall p \in \text{Domain}(\mu) : p < \nu(i)) \wedge \text{Injective}(\mu)\}$
$\gamma(\triangleright_i)$	$= \{(\mu, \nu) \mid (\forall p \in \text{Domain}(\mu) : p > \nu(i)) \wedge \text{Injective}(\mu)\}$
$\gamma(e \cdot e')$	$= \gamma(e) \cap \gamma(e')$
$\gamma(e \parallel e')$	$= \{(\mu, \nu) \mid \exists \mu_e, \mu_{e'} : (\mu_e, \nu) \in \gamma(e) \wedge (\mu_{e'}, \nu) \in \gamma(e') \wedge \text{Range}(\mu_e) \cap \text{Range}(\mu_{e'}) = \emptyset$ $\wedge (\forall p \in \text{Domain}(\mu) : \mu(p) = \mu_e(p) \vee \mu(p) = \mu_{e'}(p))\}$
$\gamma(e \# e')$	$= \{(\mu, \nu) \mid \exists \mu_e, \mu_{e'} : (\mu_e, \nu) \in \gamma(e) \wedge (\mu_{e'}, \nu) \in \gamma(e')$ $\wedge (\forall p \in \text{Domain}(\mu) : \mu(p) = \mu_e(p) \vee \mu(p) = \mu_{e'}(p))\}$
$\gamma(e^n)$	$= \{(\mu, \nu) \mid \exists \mu_e : (\mu_e, \nu) \in \gamma(e) \wedge \text{Domain}(\mu) = \text{Domain}(\mu_e)\}$

Table 3: Concretization functions for direction expressions

To ensure  $\alpha$  and  $\gamma$  preserve the relative precision over sets of connections and directions. We prove that  $\alpha$  and  $\gamma$  are monotonic over  $\sqsubseteq$ :

$$\forall C \subseteq \text{Conn} : C \subseteq C' \implies \alpha(C) \sqsubseteq \alpha(C') \quad \forall d \in \text{Dir} : d \sqsubseteq d' \implies \gamma(d) \subseteq \gamma(d') \quad (1)$$

In addition, the following equation says that  $\gamma$  is a safe approximation and that  $\gamma$  is the inverse image of  $\alpha$ , ensuring that no information is lost during concretization:

$$\forall C \subseteq \text{Conn} : C \subseteq \gamma(\alpha(C)) \quad \forall d \in \text{Dir} : d = \alpha(\gamma(d)) \quad (2)$$

Proof of (2) can be found in Section B.5. The following result summarizes these properties in the unified model of Cousot and Cousot [5] for program analysis.

**Theorem 1** *Abstraction and concretization functions  $\alpha$  and  $\gamma$  define a Galois connection from the lattice of sets of connections to the lattice of directions.*

We use this result in the appendix to prove properties of the direction abstraction (e.g., algebraic properties of operators), and to prove the analysis safety (i.e., it computes a correct approximation for each container layout).

## 4 Points-to Graphs

### 4.1 The Node Naming Scheme

In points-to graphs, heap objects are represented by *heap nodes* named as  $n_X^H$ :  $X$ , the *reference set*, holds variables that are *simultaneously* pointing to the objects at this moment; and  $H$ , the *history set*, records all variables that *ever pointed* to the objects at some past execution. Among heap nodes, we further distinguish between *object nodes* and *summary nodes*: an object node abstracts a *single object* and has a *non-empty* reference set; a summary node abstracts a *pool of objects* and has an *empty* reference set. Object and summary nodes are distinguished by predicates  $\text{Obj}$  and  $\text{Sum}$ , respectively. Intuitively, object nodes keep precise information about single objects, thus, enable “killing” on assignments to object fields. Summary nodes, on the other hand, provide a mechanism to group objects together, which is necessary to make the graph finite. Additional remarks about this naming scheme are,

- The node naming scheme using reference sets was first introduced by [18]. This is used to determine the *compatibility* of heap nodes: two nodes are *compatible* if they may *simultaneously* abstract (non-empty sets of) heap object. Since any variable may not simultaneously point to multiple heap objects, nodes  $n_X^H$  and  $n_{X'}^H$ , with *different* but *overlapping* reference sets (i.e.,  $X \neq X' \wedge X \cap X' \neq \emptyset$ ) can *not*

<sup>3</sup>To simplify the formalism, heap objects and container positions are labeled by integers.

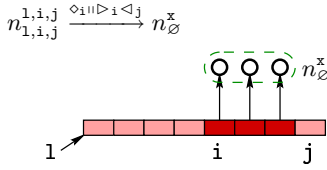


Figure 4: Container edges and direction labels

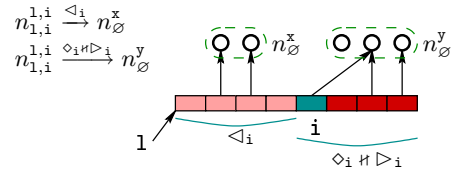


Figure 5: Multiple edges

*simultaneously* abstract heap objects. In other words, they are not compatible. In a similar way,  $n_X^H$  and  $n_{X'}^{H'}$ , with  $X = X' \neq \emptyset$  and  $H \neq H'$  are *not* compatible either. Node compatibility can improve the analysis' precision because edges that connect incompatible nodes are spurious and can be removed from the graph. We introduce predicate `Compat` to compute such information:

$$\text{Compat}(n_X^H, n_{X'}^{H'}) \stackrel{\text{def}}{\iff} (X = X' \wedge H = H') \vee X \cap X' = \emptyset.$$

- History sets make our node naming scheme different from the original one proposed in [19] and further extentions proposed in [18, 20, 4]. Indeed, history sets are used to create more object/summary nodes, preserving useful context information. This is similar to other heuristics [18, 20, 4], which resorts to allocation sites or data types to achieve the same goal.
- The points-to graph describes a *partition of memory*. In fact, for any given runtime program point, a heap object can only be abstracted by *one* heap node in the graph. However, the node naming scheme is *store-less*, meaning that nodes in a points-to graph do *not* describe a *fixed* partition of memory/store. In other words, whether an object is abstracted by a node or another depends on the program point. This store-less approach is critical for the precise handling of random container updates.

## 4.2 Definition of Points-to Graphs

We define a points-to graph  $G = (V, N, R, E, C)$  as follows.

- $V$  is the set of *program variables*.
- $N$  is the set of *heap nodes*.
- $R \subseteq V \times \{o \in N \mid \text{Obj}(o)\}$  is the set of *variable edges*. It holds directed edges of the form  $v \rightarrow o$ , indicating that variable  $v$  *must* point to the *single* object abstracted by *object* node  $o$ .
- $E \subseteq N \times \text{Fld} \times N$  is the set of *field edges*. It holds labeled directed edges of the form,  $o \xrightarrow{f} o'$ , indicating that *field*  $f$  of an object abstracted by  $o$  *may* point to an object abstracted by  $o'$ .
- $C \subseteq N \times \text{Dir} \times N$  is the set of *container edges*. It holds labeled directed edges of the form,  $l \xrightarrow{d} o$ , indicating that *every* connection between a container abstracted by  $l$  and an object abstracted by  $o$  satisfies  $d$ .

Note that, edges in  $R$  and  $E$  represent physical pointers, whereas  $C$  only abstracts logical connections between containers and their elements. Figure 4 and 5 give examples of container layouts and the corresponding container edges.<sup>4</sup>

To re-address the node compatibility and the partition property, Figure 6 gives two examples where  $x$  and  $y$  may point to the same location or not. The two heap layouts satisfy the partition property: either node  $n_{x,y}^{x,y}$  abstracts an object simultaneously referenced by  $x$  and  $y$  or nodes  $n_x^{x,y}$  and  $n_y^y$  abstract distinct objects.

To ensure the graph is finite, two important steps are applied after each points-to graph update:

- A garbage-collection scheme removes any node (along with its adjacent edges) when it is not reachable from any variable node. Container edges labeled  $\emptyset$  are removed as well.

<sup>4</sup>A container node can have multiple outgoing edges (targeting distinct nodes). Notice an object abstracted by  $n_\emptyset^y$  is not required to be an element of the container.

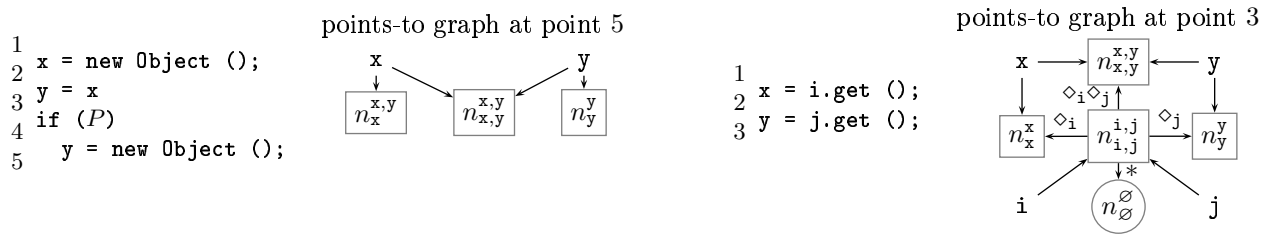


Figure 6: Node compatibility

- To ensure that a finite set of iterator variables generate a *finite number* of directions, labels are simplified according to equivalence of direction expressions. The formal rules are given in the Appendix. A special case is that directions of the form  $e^ne$  are replaced by  $e^n$ . This step can be seen as a “widening” transformation, which conservatively drops injective properties when they are unlikely to bring any precision improvement. Local safety proofs in Section B.6.

## 5 Points-to Analysis

The points-to analysis to be presented is intra-procedural except for the dedicated interpretation of container-related methods.

### 5.1 An Intuitive Flavor of the Analysis

Let us first illustrate the points-to analysis in the following example that reverses a list. We would like to infer that the `reverse()` method preserves combness.

```

1 static void reverse (Container c, Container d) {
2   Iterator i = c.last ();
3   Iterator j = d.first ();
4   while (!i.atBegin ()) {
5     Object o = i.get ();
6     j.put (o);
7     j.advance ();
8     i.retreat ();
9   }
}

```

Figure 7 describes the first two iterations of the `while` loop. For each run-time program point, it gives a concrete heap layout in parallel with the points-to graph corresponding to the layout. Initially, we assume that `c` and `d` point to distinct containers, that `c` is a comb, and that `d` is empty. In the beginning, the graph holds three nodes: object nodes  $n_c^c$  and  $n_d^d$  which abstract container objects referenced by `c` and `d`, and a summary node,  $n_\emptyset^\emptyset$ , which abstracts elements of  $n_c^c$ . The container edge between  $n_c^c$  and  $n_\emptyset^\emptyset$  is labeled  $*$ , describing an injective connection.

**Point 4<sub>1</sub>.** The `while` loop is entered for the first time. Iterator `i` (resp. `j`) is attached to the container referenced by `c` (resp. `d`). These container nodes are renamed to  $n_{c,i}^{c,i}$  and  $n_{d,j}^{d,j}$  to reflect the attachments of the iterators.

**Point 5<sub>1</sub>,** `o = i.get()`. The object attached at the position of `i` is now assigned to variable `o`. Since iterator `i` is associated with  $n_{c,i}^{c,i}$ , this object must be abstracted by  $n_\emptyset^\emptyset$  (elements of  $n_{c,i}^{c,i}$ ). Then, the object is “materialized” [18] from summary node  $n_\emptyset^\emptyset$  into a new object node. This new object node, now pointed to by `o`, is named  $n_o^o$ . The connection from  $n_{c,i}^{c,i}$  to  $n_o^o$  is captured by  $\diamond_i$ .

**Point 6<sub>1</sub>,** `j.put(o)`. Since iterator `j` is associated with container node  $n_{d,j}^{d,j}$ ,  $n_o^o$  is connected to  $n_{d,j}^{d,j}$  by a container edge labeled  $\diamond_j$ .

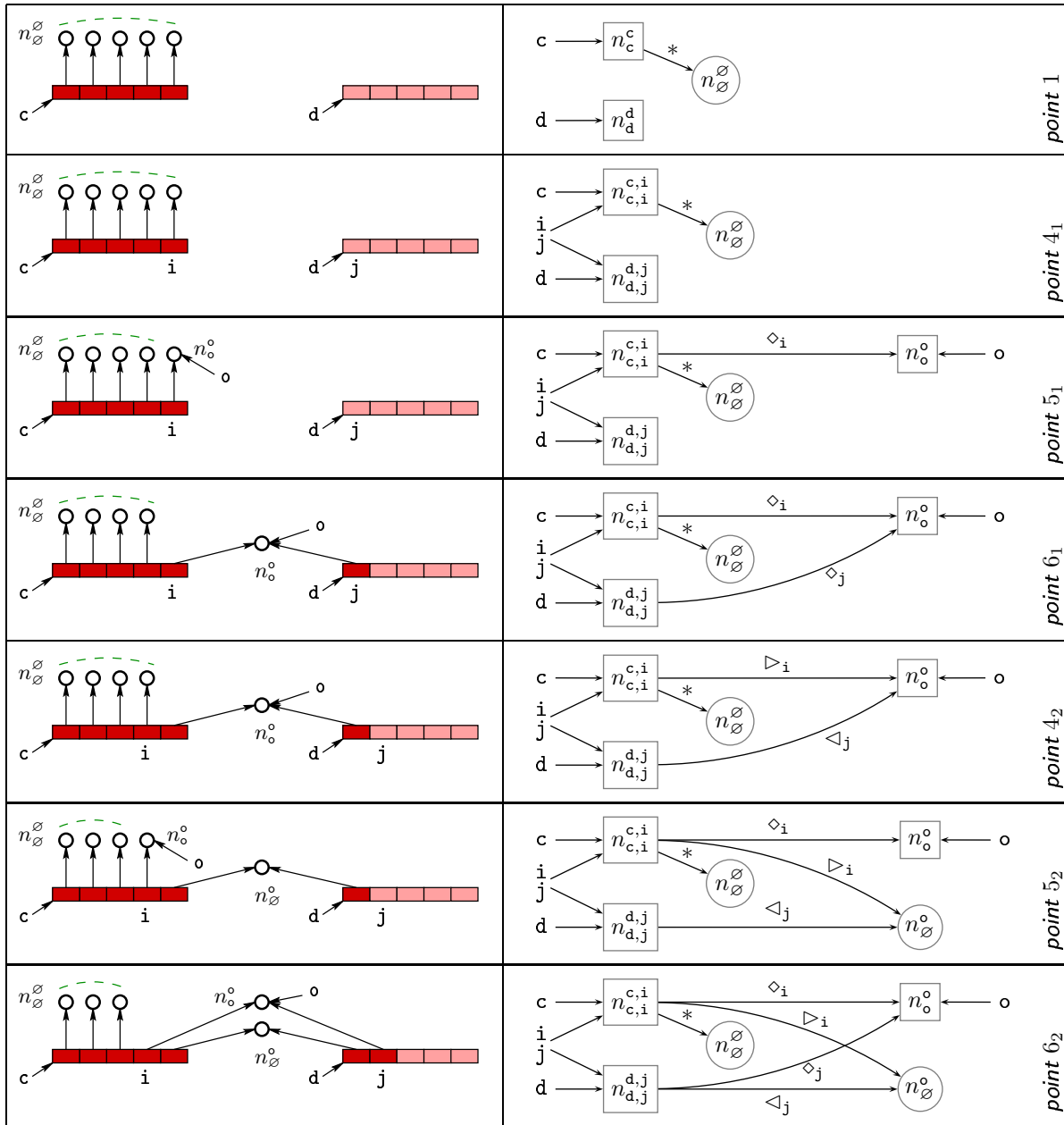


Figure 7: Points-to analysis for reverse

**Point 4<sub>2</sub>.** Iterator  $i$  (resp.  $j$ ) is moved backward (resp. forward) before entering the next iteration. Accordingly, the direction labels are updated:  $\diamond_i$  is replaced by  $\triangleright_i$ ,  $\diamond_j$  is replaced by  $\triangleleft_j$ .

**Point 5<sub>2</sub>,**  $o = i.get()$ . First, we “kill” the outgoing edges of  $o$ , i.e., removing the edge from  $o$  to  $n_o^{\circ}$ . This also involves removing  $o$  from the reference set of  $n_o^{\circ}$ . As a result,  $n_o^{\circ}$  becomes a summary node  $n_{\emptyset}^{\circ}$ . We call such a process *node summarization*. Then, the object attached at the position of  $i$  needs to be extracted from the elements of  $n_{c,i}^{\circ}$ . Note that, both  $n_{\emptyset}$  and  $n_o^{\circ}$  are attached at  $n_{c,i}^{\circ}$ . However, since the edge targetting  $n_{\emptyset}$  is labeled as  $\triangleright_i$ , which is disjoint with  $\diamond_i$ , the object can not be abstracted by  $n_o^{\circ}$ . Therefore, we only materialize from  $n_{\emptyset}$ . Again, the new object node is named  $n_{\emptyset}^{\circ}$ .

**Point 6<sub>2</sub>**, `j.put(o)`. the same as **Point 6<sub>1</sub>**. Note that, the outgoing edges of  $n_{d,j}^{d,j}$  and  $n_{c,i}^{c,i}$  are labeled by injective directions, meaning that both containers are comb.

Further iterations follow a similar pattern. A fixed point is actually reached at the third iteration. To conclude this example, we would like to emphasize three key properties of the analysis that make us successfully detect the preservation of combness.

- Node materialization helps to restore precise points-to information when an object abstracted by a summary node is assigned to a variable (e.g., `o = i.get()`).
- Because of history sets, after the “kill” of variable `o`,  $n_o^o$ —which abstracts the “tail” elements of  $n_{c,i}^{c,i}$ —is summarized into  $n_{\emptyset}^o$ , avoiding spurious confusion with “head” elements abstracted by  $n_{\emptyset}^o$ .
- Comparing  $\diamond_i$  and  $\triangleright_i$  through `Disjoint( $\diamond_i, \triangleright_i$ )` avoids materialization from  $n_{\emptyset}^o$ .

## 5.2 Strong Nullification of Variables

A strong nullification, for a points-to graph, means the “kill” (i.e., removal) of points-to edges. More specifically, a strong nullification of variables/fields/containers correspond to the removal of variable/field/container edges, respectively. The ability to “kill” is key to the high precision of this analysis. We define an auxiliary “kill” function to implement the strong nullification of *variables*. The “kill” functions for fields or container positions are implemented directly into transfer functions.

Given a points-to graph,  $G = (V, N, R, E, C)$ , `KillVariablev(G)` implements the strong nullification of  $v$ , i.e.,  $v = \text{null}$ . The function first removes all outgoing edges of  $v$ . Then,  $v$  is removed from the reference set of each node, i.e.,  $n_X^H$  is renamed to  $n_{X \setminus \{v\}}^H$ . If node  $n_{X \setminus \{v\}}^H$  already exists in the graph, edges of the two nodes need to be merged to ensure *at most* one edge between any two nodes in the graph.

- When both have an incoming container edge from the *same* container node, merging of the two edges shall preserve combness. We define operator  $\cup^{\text{in}}$  to implement such a merge:

$$C \cup^{\text{in}} C' \stackrel{\text{def}}{=} (C \setminus C') \cup (C' \setminus C) \cup \{o \xrightarrow{dnd'} o' \mid o \xrightarrow{d} o' \in C \wedge o \xrightarrow{d'} o' \in C'\}.$$

- When both merging nodes have an outgoing container edge to the *same* node, combness is not preserved because a single element may be attached in multiple containers. Operator  $\cup^{\text{out}}$  implements this merge:

$$C \cup^{\text{out}} C' \stackrel{\text{def}}{=} (C \setminus C') \cup (C' \setminus C) \cup \{o \xrightarrow{dnd'} o' \mid o \xrightarrow{d} o' \in C \wedge o \xrightarrow{d'} o' \in C'\}.$$

In the points-to analysis, `KillVariable` is the *only way* to generate a summary node. `KillVariablev(G)` is defined below.<sup>5</sup> Examples are given in Figure 8.

$$\begin{aligned} \text{KillVariable}_v(V, N, R, E, C) &= (V, N', R', E', C') \\ (N', R', E', C') &\leftarrow (N, R, E, C) \\ n_X^H \in N \wedge v \in X &: N' \leftarrow N' \setminus \{n_X^H\} \cup \{n_{X \setminus \{v\}}^H\}; R' \leftarrow R' \setminus \{v \rightarrow n_X^H\} \\ o \xrightarrow{f} n_X^H \in E \wedge v \in X &: E' \leftarrow E' \setminus \{o \xrightarrow{f} n_X^H\} \cup \{o \xrightarrow{f} n_{X \setminus \{v\}}^H\} \\ n_X^H \xrightarrow{f} o \in E \wedge v \in X &: E' \leftarrow E' \setminus \{n_X^H \xrightarrow{f} o\} \cup \{n_{X \setminus \{v\}}^H \xrightarrow{f} o\} \\ o \xrightarrow{d} n_X^H \in C \wedge v \in X &: C' \leftarrow C' \setminus \{o \xrightarrow{d} n_X^H\} \cup^{\text{in}} \{o \xrightarrow{d} n_{X \setminus \{v\}}^H\} \\ n_X^H \xrightarrow{d} o \in C \wedge v \in X \wedge \text{Compat}(n_X^H, n_{X \setminus \{v\}}^H) &: C' \leftarrow C' \setminus \{n_X^H \xrightarrow{d} o\} \cup^{\text{out}} \{n_{X \setminus \{v\}}^H \xrightarrow{d} o\} \\ n_X^H \xrightarrow{d} o \in C \wedge v \in X \wedge \neg \text{Compat}(n_X^H, n_{X \setminus \{v\}}^H) &: C' \leftarrow C' \setminus \{n_X^H \xrightarrow{d} o\} \cup^{\text{in}} \{n_{X \setminus \{v\}}^H \xrightarrow{d} o\}. \end{aligned}$$

<sup>5</sup>For this function and graph transformers to come, rules of the form  $C : A$  use a side-effect notation where  $C$  is a “for-all” condition which “triggers” action  $A$ . Rules apply in-order but conditions are exclusive when present.

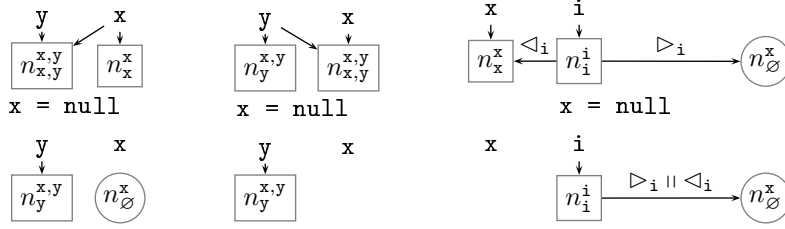


Figure 8: Strong nullifications of variables

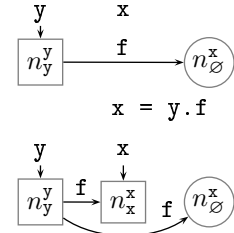


Figure 9: Materialization

### 5.3 Node Materialization

For best precision, an object abstracted by a summary node needs to be “materialized” to an object node when it is assigned to a variable [18]. Figure 9 gives an example of such a process. As a result, in a points-to graph, variables can not directly point to summary nodes. We define an auxiliary function to implement such a graph transformation:  $\text{Materialize}_{v,n_X^H}(G)$  “extracts” an object node from a summary node  $n_{\emptyset}^H$  and connects it with variable  $v$ . Basically, the function duplicates  $n_{\emptyset}^H$  (including its adjacent edges) with a name  $n_v^{H \cup \{v\}}$  and adds  $v \rightarrow n_v^{H \cup \{v\}}$  to the graph.

We also materialize nodes through container edges. When materializing an object node for  $i.\text{get}()$ , one needs a new container edge to record the fact that the materialized node is attached at  $\diamond_i$ . We thus generalize the concept of materialization to “transfer” an object from an *object or summary* node to an object node.  $\text{Materialize}_{v,n_X^H}(G)$  duplicates  $n_X^H$  (including its adjacent edges) with a name  $n_{X \cup \{v\}}^{H \cup \{v\}}$  and links  $v$  to the duplicated node.<sup>6</sup> Materialize is formally defined below.

$$\begin{aligned} \text{Materialize}_{v,n_X^H}(V, N, R, E, C) &= (V, N \cup \{n_{X \cup \{v\}}^{H \cup \{v\}}\}, R \cup \{v \rightarrow n_{X \cup \{v\}}^{H \cup \{v\}}\} \cup \{w \rightarrow n_{X \cup \{v\}}^{H \cup \{v\}} \mid w \rightarrow n_X^H \in R\}, \\ E \cup \{o \xrightarrow{f} n_{X \cup \{v\}}^{H \cup \{v\}} \mid o \xrightarrow{f} n_X^H \in E\} \cup \{n_{X \cup \{v\}}^{H \cup \{v\}} \xrightarrow{f} o \mid n_X^H \xrightarrow{f} o \in E\}, \\ C \cup \{o \xrightarrow{d} n_{X \cup \{v\}}^{H \cup \{v\}} \mid o \xrightarrow{d} n_X^H \in C\} \cup \{n_{X \cup \{v\}}^{H \cup \{v\}} \xrightarrow{d} o \mid n_X^H \xrightarrow{d} o \in C\}). \end{aligned}$$

### 5.4 Transfer Functions

In the following,  $X[b/a]$  denotes the substitution of  $a$  with  $b$  in each element of set  $X$ .

#### 5.4.1 The Classical Part

This part discusses transfer functions for pointer assignments that are not specific to containers. We use an auxiliary function  $\text{Link}_{v,n_X^H}(G)$  to connect a variable  $v$  to an *object* node  $n_X^H$  ( $X \neq \emptyset$ ) and then to rename  $n_X^H$  to  $n_{X \cup \{v\}}^{H \cup \{v\}}$ . Link assumes that a strong nullification of variable  $v$  is applied before, therefore  $n_{X \cup \{v\}}^{H \cup \{v\}}$  does not name-conflict with any other node. Here is the formal definition of Link:

$$\text{Link}_{v,n_X^H}(V, N, R, E, C) = (V, N[n_{X \cup \{v\}}^{H \cup \{v\}}/n_X^H], R \cup \{v \rightarrow n_{X \cup \{v\}}^{H \cup \{v\}}\}, E[n_{X \cup \{v\}}^{H \cup \{v\}}/n_X^H], C[n_{X \cup \{v\}}^{H \cup \{v\}}/n_X^H]).$$

Transfer functions are specified in Figure 10.

$v = w$ ,  $v = \text{null}$ , and  $v = \text{new}$ . It involves a strong nullification of  $v$ , followed by linking  $v$  to an object node.

$v.f = w$ . It implements a strong nullification of object field  $v.f$ : it removes every outgoing field-edge of *object* nodes pointed by  $v$  that are labeled by  $f$ . As opposed to KillVariable, this involves no renaming, thus no merging of nodes is needed.

<sup>6</sup>Safety of materialization from an object node (when  $X$  is not empty) is similar to safety of joining points-to graphs at control-flow merges in Sagiv, Reps and Wilhelm technique [18]. Indeed, nodes  $n_X^H$  and  $n_{X \cup \{v\}}^{H \cup \{v\}}$  may not simultaneously abstract a concrete object since they are not *compatible*.



$v = \text{null}$	$G \leftarrow \text{KillVariable}_v(G)$
$v = \text{new Object}()$	$G \leftarrow \text{KillVariable}_v(G)$ $N \leftarrow N \cup \{n_v^v\}; R \leftarrow R \cup \{v \rightarrow n_v^v\}$
$v = w$	$G \leftarrow \text{KillVariable}_v(G)$ $w \rightarrow n_X^H \in R : G \leftarrow \text{Link}_{v, n_X^H}(G)$
$v = w.f$	$G \leftarrow \text{KillVariable}_v(G)$ $w \rightarrow o \in R \wedge o \xrightarrow{f} n_X^H \in E \wedge \text{Obj}(n_X^H) : G \leftarrow \text{Link}_{v, n_X^H}(G)$ $w \rightarrow o \in R \wedge o \xrightarrow{f} n_\emptyset^H \in E : G \leftarrow \text{Materialize}_{v, n_\emptyset^H}(G)$
$v.f = w$	$v \rightarrow o \in R \wedge o \xrightarrow{f} o' \in E : E \leftarrow E \setminus \{o \xrightarrow{f} o'\}$ $v \rightarrow o \in R \wedge w \rightarrow o' \in R \wedge \text{Compat}(o, o') : E \leftarrow E \cup \{o \xrightarrow{f} o'\}$

Figure 10: Transfer functions of non-container specific assignments

$i = l.\text{first}()$	$G \leftarrow \text{KillVariable}_i(G)$
$i = l.\text{last}()$	$C \leftarrow C[*/\diamond_i, */\triangleleft_i, */\triangleright_i]$ $l \rightarrow n_X^H \in R : G \leftarrow \text{Link}_{i, n_X^H}(G)$
$i.\text{advance}()$	$C \leftarrow C[\triangleleft_i/\diamond_i, */\triangleright_i]$
$i.\text{retreat}()$	$C \leftarrow C[\triangleright_i/\diamond_i, */\triangleleft_i]$

Figure 11: Transfer functions of iterator operations

$v = w.f$ . It starts with a strong nullification of  $v$ . Then, if  $w.f$  points to an object node,  $v$  is directly linked to that object node. On the other hand, if  $w.f$  points to a summary node, an object node is first materialized, then it is linked to  $v$ .

#### 5.4.2 Iterator Operations

We only consider iterator variables that are *unaliased*. This simplifying assumption has two consequences:

- Calling a method on an iterator variable will not affect iterator objects referenced by other variables.
- The attachment of an iterator to a container node can be represented as a *variable edge* with the iterator variable appearing in the reference and history sets of the container node. Such variable edges capture *must-attachment* properties,<sup>7</sup> which is key to the strong nullification of container fields.

Figure 11 gives transfer functions for iterator operations that involve no container update/access. The explanation is given below:

$i=l.\text{first}()$  and  $i=l.\text{last}()$ . It involves a strong nullification of iterator  $i$  and linking  $i$  to container  $l$ .

In addition, since  $i$  is re-attached, previous direction labels involving  $i$  become obsolete, thus, replaced by  $*$ .<sup>8</sup>

$i.\text{advance}()$ . Since  $i$  moves forward, previous direction labels involving  $i$  are updated:  $\diamond_i$  is replaced by  $\triangleleft_i$  and  $\triangleright_i$  by  $*$ .

$i.\text{retreat}()$ . Since  $i$  moves backward, the previous direction labels involving  $i$  are updated:  $\diamond_i$  is replaced by  $\triangleright_i$  and  $\triangleleft_i$  by  $*$ .

<sup>7</sup>Only *may-point-to* properties would be available if the attachment was indirect, i.e., if we connected an iterator *object node* to the container.

<sup>8</sup>As a consequence, a direction label of a container edge may only involve iterators currently attached to the container. Moreover, any outgoing container edge from summary container node has direction label  $*$  or  $*'$ .

### 5.4.3 Container Operations

This section discusses iterator operations that update or access container elements. Different semantics have been proposed for such operators. We choose one that is both realistic and easy to adapt to a particular library: none of the iterator operations implicitly move iterators. Specifically, `i.delete()` removes the object at the position of `i` leaving `i` “between” two legal positions;<sup>9</sup> `i.insert(v)` adds an object *after* the position of `i`; `i.put(v)` and `i.get()` occur at the position of `i`.

<code>i.delete()</code>	$i \rightarrow l \in R \wedge l \xrightarrow{d} o \in C : C \leftarrow C \setminus \{l \xrightarrow{d} o\} \cup \{l \xrightarrow{d \cdot (\triangleleft_i \parallel \triangleright_i)^n} o\}$
<code>i.insert(v)</code>	$i \rightarrow l \in R \wedge v \rightarrow o \in R \wedge \text{Compat}(l, o) : C \leftarrow C \cup^n \{l \xrightarrow{\triangleright_i} o\}$
<code>i.put(v)</code>	$i \rightarrow l \in R \wedge v \rightarrow o \in R \wedge l \xrightarrow{d} o' \in C \wedge o \neq o' : C \leftarrow C \setminus \{l \xrightarrow{d} o'\} \cup \{l \xrightarrow{d \cdot (\triangleleft_i \parallel \triangleright_i)^n} o'\}$ $i \rightarrow l \in R \wedge v \rightarrow o \in R \wedge \neg(l \xrightarrow{d} o \in C) \wedge \text{Compat}(l, o) : C \leftarrow C \cup \{l \xrightarrow{\diamond_i} o\}$ $i \rightarrow l \in R \wedge v \rightarrow o \in R \wedge l \xrightarrow{d} o \in C \wedge \neg(\diamond_i \in d) : C \leftarrow C \setminus \{l \xrightarrow{d} o\} \cup \{l \xrightarrow{d \circ \diamond_i} o\}$ $i \rightarrow l \in R \wedge v \rightarrow o \in R \wedge l \xrightarrow{d} o \in C \wedge \diamond_i \in d : C \leftarrow C \setminus \{l \xrightarrow{d} o\} \cup \{l \xrightarrow{d \circ \diamond_i} o\}$
<code>v = i.get()</code>	$G \leftarrow \text{KillVariable}_v(G)$ $i \rightarrow l \in R \wedge l \xrightarrow{d} n_X^H \in C \wedge \text{Obj}(n_X^H) \wedge \text{Coincide}(\diamond_i, d) : G \leftarrow \text{Link}_{v, n_X^H}(G)$ $i \rightarrow l \in R \wedge l \xrightarrow{d} n_X^H \in C \wedge \neg(\text{Obj}(n_X^H) \wedge \text{Coincide}(\diamond_i, d)) \wedge \neg \text{Disjoint}(\diamond_i, d) \wedge \text{Comb}(d) :$ $G \leftarrow \text{Materialize}_{v, n_X^H}(G); C \leftarrow C \setminus \{l \xrightarrow{d} n_{X \cup \{v\}}^{H \cup \{v\}}\} \cup \{l \xrightarrow{\diamond_i \cdot d} n_{X \cup \{v\}}^{H \cup \{v\}}\}$ $i \rightarrow l \in R \wedge l \xrightarrow{d} n_X^H \in C \wedge \neg(\text{Obj}(n_X^H) \wedge \text{Coincide}(\diamond_i, d)) \wedge \neg \text{Disjoint}(\diamond_i, d) \wedge \neg \text{Comb}(d) :$ $G \leftarrow \text{Materialize}_{v, n_X^H}(G)$

Figure 12: Transfer functions of container operations

Transfer functions are given in Figure 12. Some examples are shown in Figure 13.

`i.delete()`. This is the container counterpart of `v.f = null`. The transfer function involves a strong nullification at the position of `i`. In other words, it removes any connection (container edge) whose domain can be described by  $\diamond_i$ . Having the property,

$$\diamond_i \cdot (\triangleleft_i \parallel \triangleright_i)^n \equiv (\diamond_i \cdot \triangleleft_i)^n \parallel (\diamond_i \cdot \triangleright_i)^n \equiv \emptyset$$

this step is implemented by a substitution of any label  $d$  by  $d \cdot (\triangleleft_i \parallel \triangleright_i)^n$ . The actual edge removal, however, occurs at the garbage collection step when the direction label of an edge is equivalent to  $\emptyset$ .

`i.insert(v)`. An object is attached immediately follows the position of `i`. The transfer function adds a container edge labeled by  $\triangleright_i$ . Operator  $\cup^n$  is required, in this case, because the new container edge may conflict with an existing one.

`i.put(v)`. This operation combines the effect of a deletion and an insertion. Let  $l$  denote any container node pointed by `i` and let  $o$  denotes any node pointed by `v`. The transfer function starts with a strong nullification at the position of `i`: all edges targetting object nodes which are *not* pointed by `v` are updated the same way as `i.delete()`.

Then, we use the notation  $\diamond_i \in d$  to indicate that  $\diamond_i$  appears in *every expression* of  $d$ .<sup>10</sup> We use an important property of the analysis.

**Lemma 1** *Let  $d$  be a direction such that  $\diamond_i \in d$ , i.e.,  $\diamond_i$  appears in every expression of  $d$ . For any configuration of iterator positions,  $d$  describes a connection whose domain is the current position of  $i$ .*

*In formal terms, for any mapping  $\nu$ , there is a mapping  $\mu_i$  such that  $\mu_{i, \nu} \in \gamma(d)$  and  $\text{Domain}(\mu_i) = \{i\}$ .*

<sup>9</sup>This is similar to the semantics of `remove()` in Java SDK

<sup>10</sup>This is equivalent to checking for  $\diamond_i$  in a flattened form expression of  $d$ , after exclusive primitive directions have been removed.

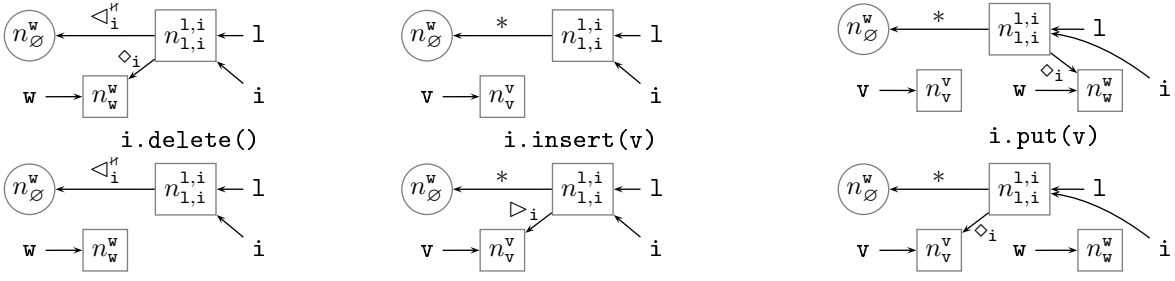


Figure 13: Container updates

There are three exclusive cases.

1. If there is no container edge targetting  $o$  from  $l$ , the transfer function simply adds  $l \xrightarrow{\diamond_i} o$  to the graph.
2. When there exists an edge  $l \xrightarrow{d} o$  such that  $\diamond_i \in d$ , the merged edge is labeled  $d_H(d \cdot \diamond_i)$ . Indeed, Lemma 1 enforces that a connection cannot satisfy  $\diamond_i$  without satisfying  $d$  as well. Therefore, the new connection from the position of  $i$  to the object abstracted by  $o$  is already described by  $d$ .
3. When there exists an edge  $l \xrightarrow{d} o$  such that  $\diamond_i \notin d$ ,  $l \xrightarrow{\diamond_i} o$  is simply merged with the edge, and label  $d$  is substituted by  $d_H \diamond_i$ . Notice direction  $d_H \diamond_i$  would be over-conservative when  $\diamond_i \in d$ : it loses the information that the new connection satisfies *both*  $\diamond_i$  and  $d$ .

$v = i.get()$ . This is the container counterpart of  $v = w.f$ . The transfer function starts with a strong nullification of variable  $v$ , followed by three exclusive cases.

1. If a node is linked to the container at the current position of  $i$  (i.e.,  $\text{Coincide}(d, \diamond_i)$ ), the node is linked to  $v$ . See examples in Figure 14.a and 15.a.
2. If a node is connected to the container through a *comb* direction that *may describe* connections satisfying  $\diamond_i$  (i.e.,  $\neg \text{Disjoint}(d, \diamond_i)$ ), we *materialize* a new object node and link  $v$  to it. Compared to  $v = w.f$ , this materialization is not limited to summary nodes. Figure 14.c gives an example of materialization from object nodes. In addition, the container edge from the container to the materialized node is labeled  $\diamond_i \cdot d$ .<sup>11</sup> See examples in Figure 14.b and 14.c.
3. This is similar to the previous case except that the container direction is *not* a comb. Materializing through a non-comb edge may create a node abstracting an object attached at several container positions: this forbids adding  $\diamond_i$  to the direction label. See examples in Figure 15.b and 15.c.

Note that, there must exist at least one container edge whose positions include  $\diamond_i$ . This means that when none of the above rules apply, a programming error is detected (getting an element from an unattached iterator).

## 5.5 Merging Edges

Eventually, we structure points-to graphs as a join semi-lattice to enable iterative data-flow analysis:

$$(V, N, R, E, C) \sqcup (V', N', R', E', C') \stackrel{\text{def}}{=} (V \cup V', N \cup N', R \cup R', E \cup E', C \cup C').$$

Preserving uniqueness of nodes and edges requires merging two edges with the same start and target nodes. This does not preserve combness in general. However, the following result shows that combness is preserved at control-flow merges during iterative analysis.

<sup>11</sup>It records that the connection between the container object and the new object node satisfies both  $d$  and  $\diamond_i$ . This is the key to proving Lemma 1.

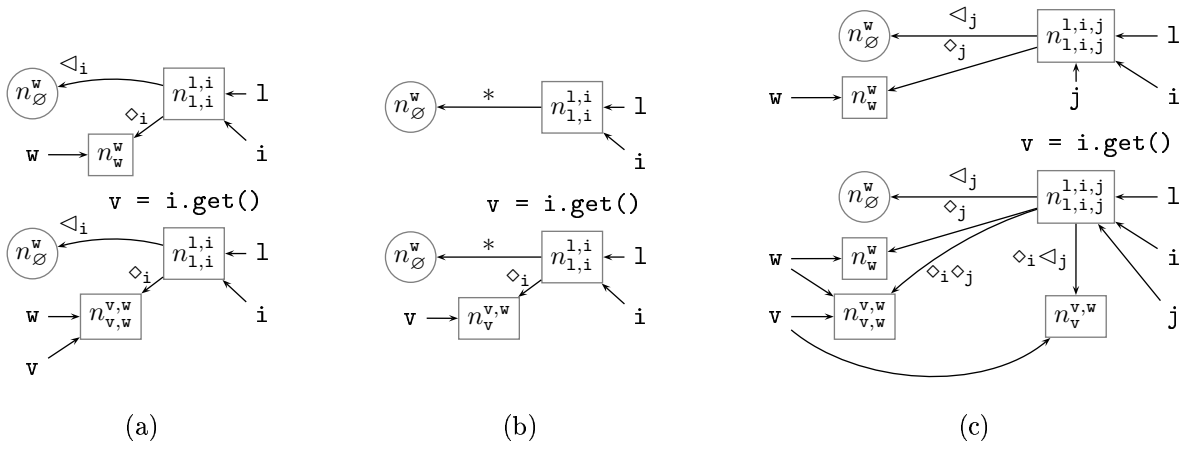


Figure 14: Retrieval through an iterator, comb case

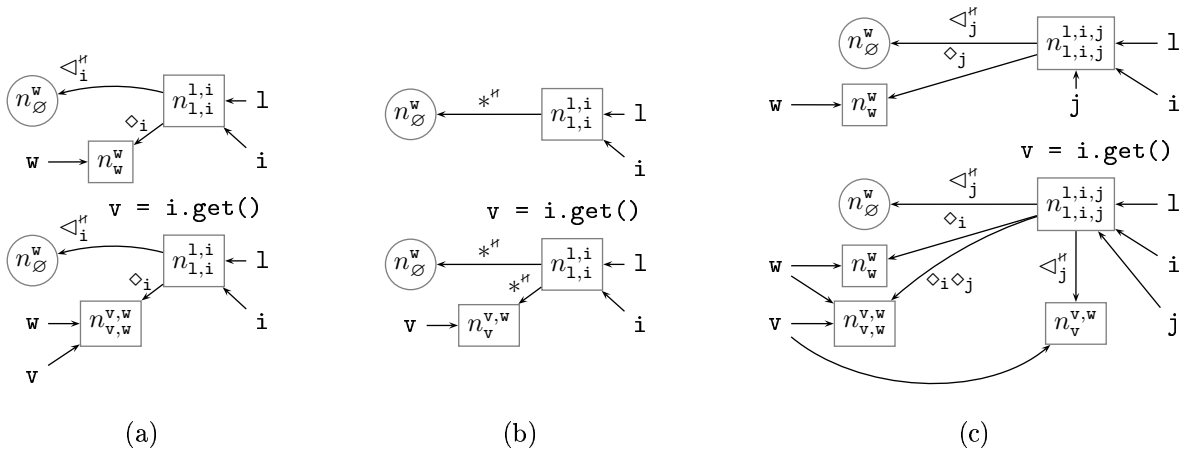


Figure 15: Retrieval through an iterator, non comb case

**Theorem 2** *When traversing a merge point in the control-flow graph, it is safe to use operator  $\sqcup$  on points-to graphs associated with incoming control-flow edges.*

Proof is given in Section B.7.

## 5.6 Correctness Issues

We already proved safe handling of directions and correctness of operator  $\sqcup$ , i.e., most of the novel features introduced in this analysis. Convergence of the iterative analysis is a consequence of points-to graphs bounded size (nodes, edges and directions) and monotonicity of transfer functions. From Theorem 1 and from local safety of each transfer function, we may apply the main result of abstract interpretation (Theorem T2 of [5] page 252) to prove global safety of our pointer analysis. However, because the current abstract interpretation does not consider points-to graph themselves, this correctness proof only deals with properties of individual container edges. Considering  $l \xrightarrow{d} o$  in a points-to graph  $G$  at program point  $p$ , and calling  $C$  the set of all possible connections at point  $p$  between containers abstracted by  $l$  and objects abstracted by  $o$ , one has  $\alpha(C) \subseteq G$ .<sup>12</sup>

<sup>12</sup>Set  $C$  can be defined as the join of concrete connections between  $l$  and  $o$  over all control-flow paths leading to  $p$  [15], a.k.a. the collecting semantics at  $p$ .

Completing the analysis’s proof is left for future work. Establishing the partition property (two nodes may not abstract the same object for a given execution) will probably be a fine point of the proof: the main idea is that new nodes are only created when (1) an object allocation occurs, yielding a distinct heap location, (2) an object node is materialized from a summary node, which implicitly states that some concrete object is “transferred” from the summary node to the object node, and (3) an object node is materialized from another object node (through method `get()`) and both are incompatible (they cannot abstract a heap object at the same time). Sagiv and al. studied a similar issue in [18] and we expect the backbone of their proof to be applicable to our case.

## 6 Illustrative Examples

Points-to analysis is now applied to small kernels of practical interest: list copy-reversal, element swapping, iterated swapping and nested traversals—a template for many (sparse) linear algebra kernels [4].

### 6.1 List Copying and Reversal

We apply the iterative algorithm to the list reversal function presented in Section 5.1. It confirms the analysis ability to discriminate current and previous positions of an iterator, using directions.

```

1  static void reverse (Container c, Container d) {
2      Iterator i = c.last ();
3      Iterator j = d.first ();
4      while (!i.atBegin ()) {
5          Object o = i.get ();
6          j.put (o);
7          j.advance ();
8          i.retreat ();
9      }
}
```

Initially, we assume that variables `c` and `d` are unaliased, that `c` is a comb—hence the comb edge with unknown direction ( $*$ ) to summary node  $n_{\emptyset}^{\circ}$ —and `d` is empty. The points-to analysis for `reverse` is sketched in Figure 16. The new container edge created at the first iteration of point 5 is labeled with  $\diamond_i$ , because materialization has been applied through a comb edge (see the transfer function for `get()`). At the second iteration of point 5, `o` does not refer an element of `d` because we could check that  $\text{Disjoint}(\diamond_i, \triangleright_i)$  and we did not materialize from  $n_{\emptyset}^{\circ}$ . The fixed point is reached after three iterations, and point 9 yields the final points-to graph: it shows that combness of `c` and `d` is preserved.

Now, if we supposed instead that `d` was connected to `c` at the beginning, the result would be quite different. Calling `reverse` on such a configuration would not preserve combness in general, and the points-to analysis detects it, as shown in Figure 17. Combness of `d` is definitively lost at point 6 because  $n_{\circ}^{\circ}$  is reachable from two distinct positions in `d`. Notice that a technique which does not allow comparison of container positions would fail to detect combness preservation: it would mix objects extracted from one container with objects stored into the other and the result would be similar to Figure 17.

Eventually, if we supposed that `c` was not a comb container at the beginning, the analysis would safely discover that `d` is not a comb either in the end. As shown in Figure 18, element retrieval still leads to materialization of an object node, but the new container edge is labeled with  $*^n$ . Then, we lose disjointness of positions accross iterations: the next element retrieval yields materializations from both  $n_{\emptyset}^{\circ}$  and  $n_{\circ}^{\circ}$ . This ends up with a container edge labeled  $\diamond_j \# \triangleleft_j$ , i.e., one element possibly being attached at several positions in container `d`.

### 6.2 Swapping Elements

Analysis is applied to a swap function; it shows the node naming scheme in action with non-primitive directions.

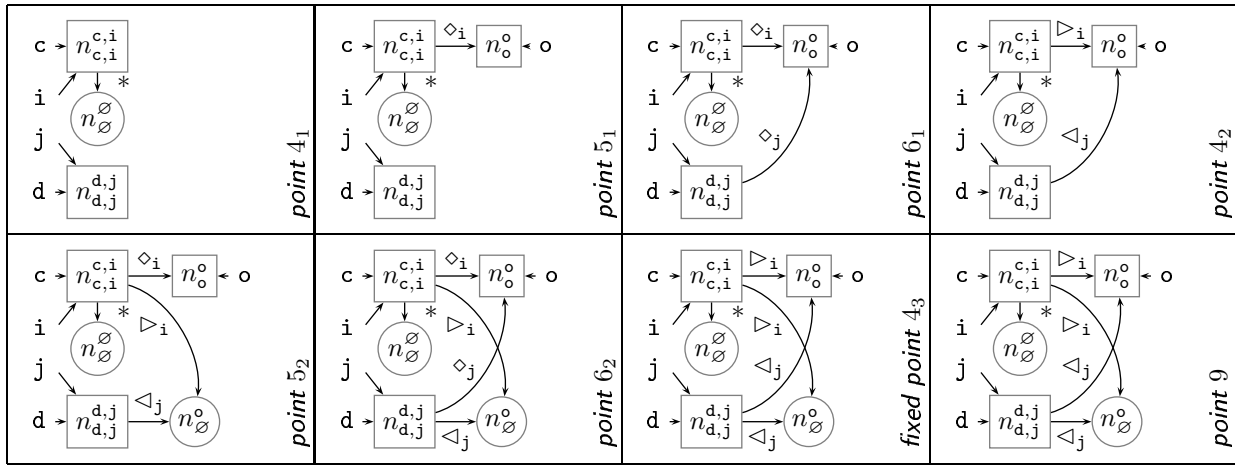


Figure 16: Points-to analysis for reverse

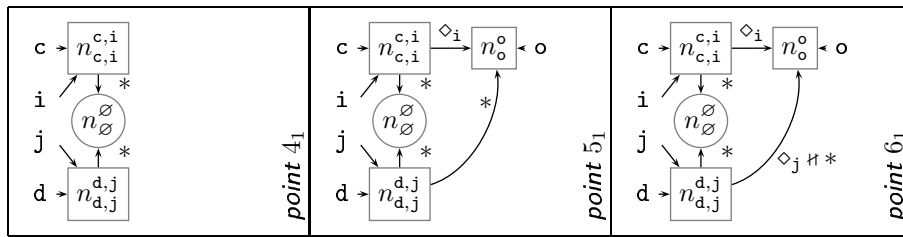


Figure 17: Points-to analysis (first steps) for reverse with connected containers

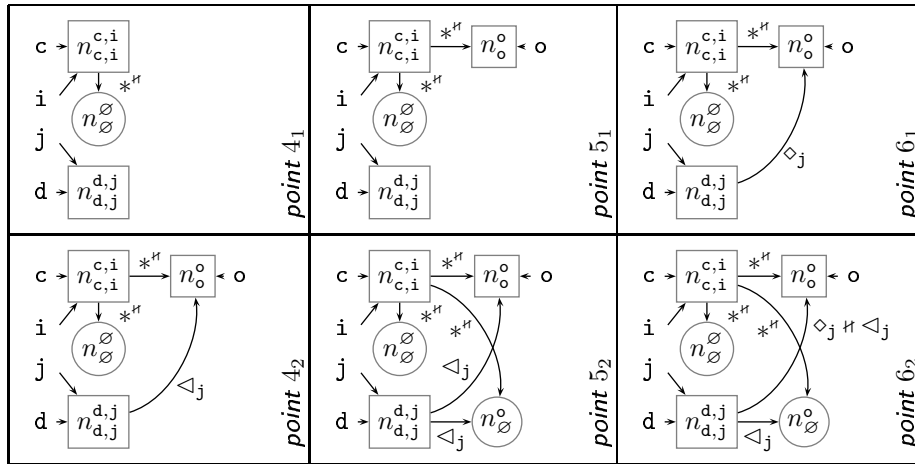


Figure 18: Points-to analysis (first steps) for reverse where container c is not a comb

```

1 static void swap (Iterator i, Iterator j) {
2     Object a = i.get ();
3     Object b = j.get ();
4     i.put (b);
5     j.put (a);
6 }

```

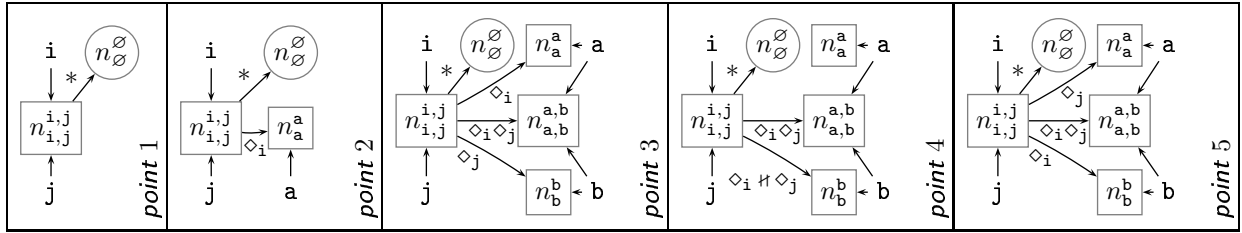


Figure 19: Points-to analysis for `swap`

We suppose that  $i$  and  $j$  are iterators attached to the *same comb container*,<sup>13</sup> but positions of  $i$  and  $j$  are arbitrary. Analysis for `swap` is sketched in Figure 19. Point 3 demonstrate materialization from both summary and object nodes. Combness of  $n_{i,j}^{i,j}$  is lost at point 4 because  $\diamond_i$  and  $\diamond_j$  cannot be compared, but is restored at point 5 thanks to the contextual knowledge captured by directions. Similarly to the destructively reversed list in [18], this example exploits flow-sensitivity and context-aware naming of object nodes.

### 6.3 In-Place Container Rotation or Reversal

Adding a loop around the previous “swap core” makes things trickier: besides evolution of iterators, an additional difficulty consists in keeping precise knowledge of possible container layouts while using summary nodes. In-place container rotation and reversal are examples of such iterated swapping schemes. They demonstrate three key analysis capabilities:

1. to handle several nodes with distinct history sets (but identical reference sets);
2. to avoid confusions between these nodes thanks to the direction abstraction;
3. and to allow iterative loss and restoration of combness.

Studying the container reversal example, which is very similar to container rotation, we automatically detect combness preservation. We are not aware of any existing technique that achieves such precision on equivalent pointer-chasing container implementations.

```

1  static void reverseInPlace (Container c) {
2      Iterator i = c.first ();
3      Iterator j = c.last ();
4      while (i.isBefore (j)) {
5          Object a = i.get ();
6          Object b = j.get ();
7          i.put (b);
8          j.put (a);
9          i.advance ();
10         j.retreat ();
11     }

```

We suppose that  $c$  is a comb container. Points-to analysis for `reverseInPlace` is sketched in Figure 20 (variable  $c$  is not represented). Combness of  $n_{i,j}^{i,j}$  is lost and restored at each iteration, but eventually preserved.

Notice that history sets are critical for the successful analysis of this example. Using the plain node naming scheme proposed in [18], or using extensions like allocation sites and type properties, we would have to merge nodes  $n_b^{a,b}$  and  $n_b^b$  at the second iteration of point 5, as well as  $n_{\emptyset}^{a,b}$  and  $n_{\emptyset}^b$  at the second iteration of

<sup>13</sup>Combness is actually the only required information: the case where iterators may *possibly* be attached to different containers can be handled with two additional object nodes  $n_i^i$  and  $n_j^j$ .

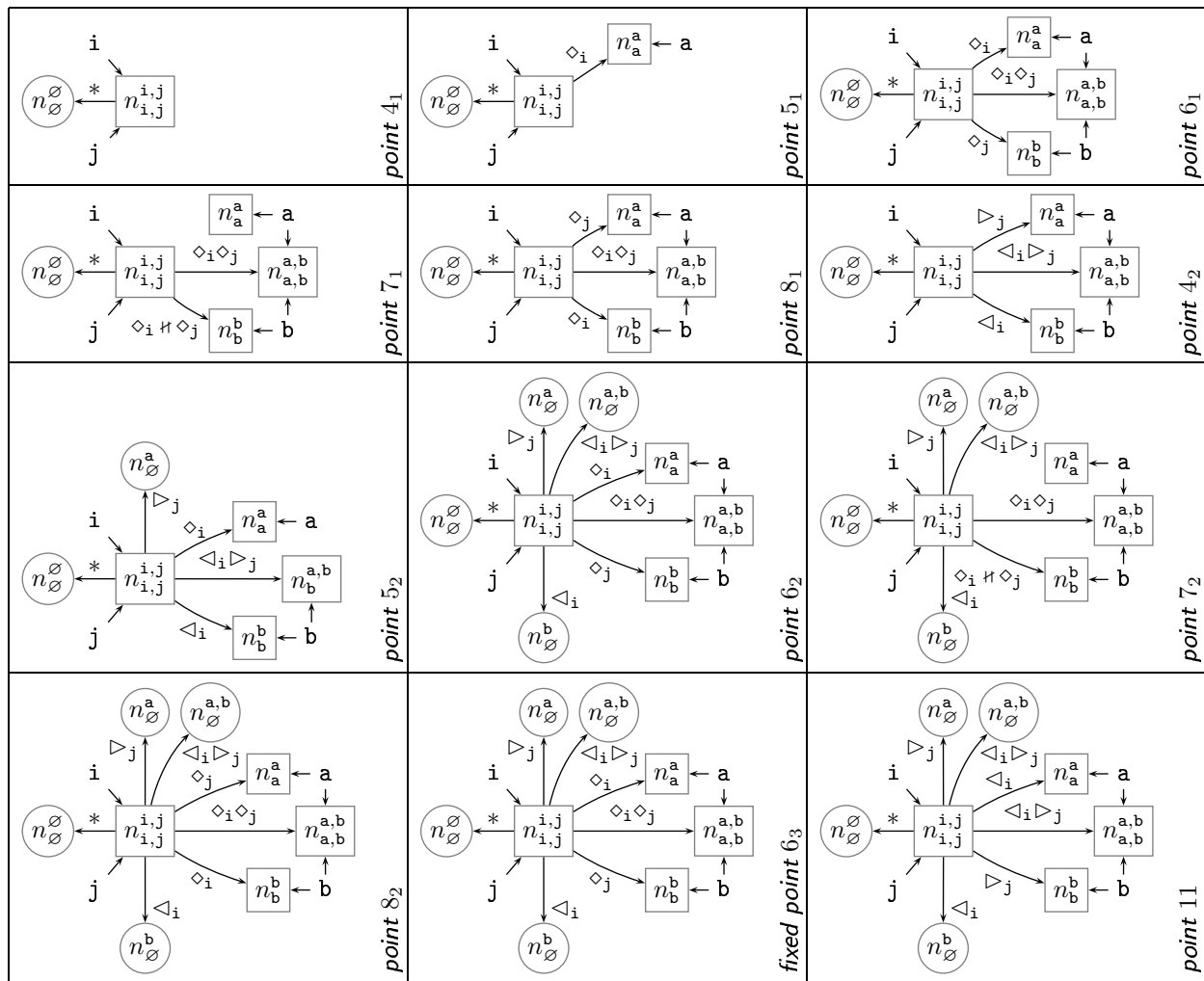


Figure 20: Points-to analysis for `reverseInPlace`

point 6. This would result in merging edges labeled with  $\triangleright_j$  and  $\triangleleft_i \triangleright_j$ , yielding  $\triangleleft_i \triangleright_j \parallel \triangleright_j \equiv \triangleleft_i$ . Therefore, `j.get()` would conservatively materialize a node  $n_{a,b}^{a,b}$  from both  $n_a^a$  and  $n_b^b$ , with an edge labeled  $(\diamond_i \parallel \triangleleft_i) \diamond_j$  instead of  $\diamond_i \diamond_j$ . This weakens the killing at the second iteration of point 7, and combness is eventually lost at point 8.

## 6.4 Nested Traversals

Multiple traversals of containers of containers seem to be a very common template for many (sparse) linear algebra kernels. We show that the analysis keeps a reasonable precision when dealing with summarized container nodes. Here is an example with two nested loops traversing a list of lists.

```

1 static void nestedTraversals (Iterator i) {
2   while (!i.atEnd ()) {
3     List l = i.get ();
4     Iterator j = l.first ();
5     while (!j.atEnd ()) {
6       Object o = j.get ();
       j.put (o);

```



```

7     j.advance ();
8     }
9     i.advance ();
10  }
11 }

```

To decide whether combness is preserved, we assume iterator  $i$  traverses a comb list of lists abstracted by object  $n_i^i$  and the two attached summary nodes  $n_\emptyset^1$  (abstraction for the second level lists) and  $n_\emptyset^\emptyset$ . Full combness of the structure is ensured by combness of edges between  $n_i^i$  and  $n_\emptyset^1$  and between  $n_\emptyset^1$  and  $n_\emptyset^\emptyset$ . The points-to analysis's main steps for `nestedTraversals` are sketched in Figure 21.

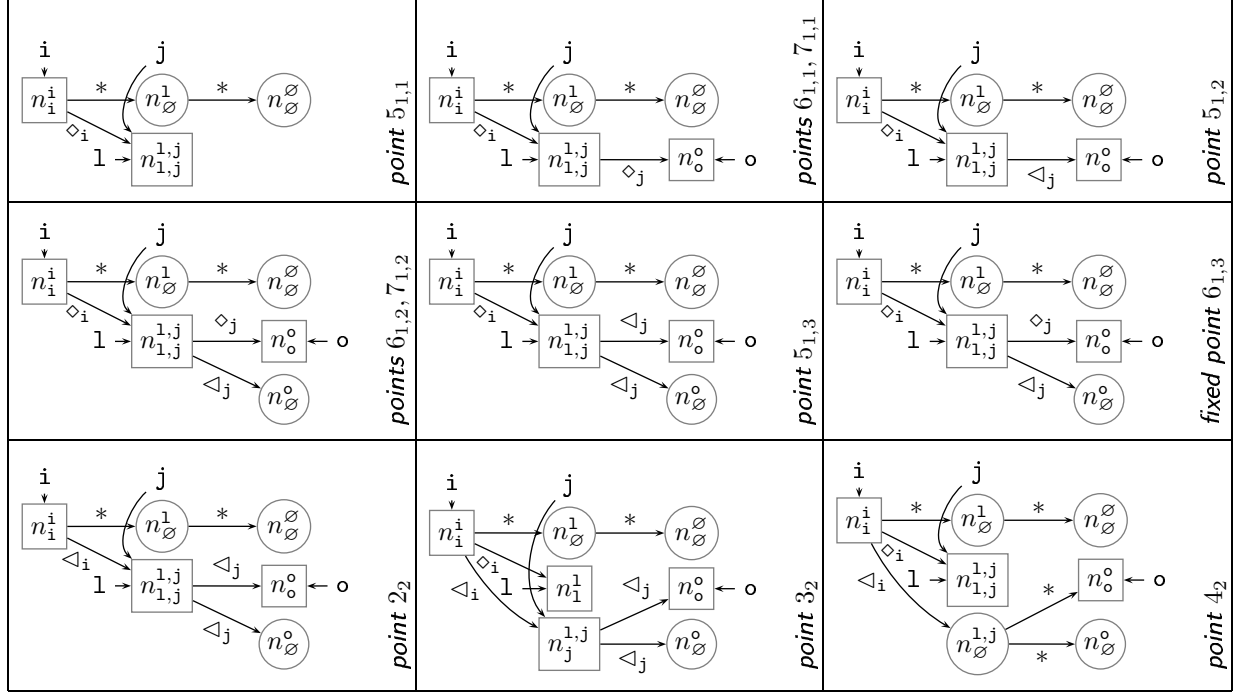


Figure 21: Points-to analysis for `nestedTraversals` (main steps)

## 7 Aliasing and Combness

Most questions related with static properties of pointers may be solved using points-to graphs.

We overload notation  $\rightarrow$  to denote the *points-to relation* itself (as opposed to points-to graph edges):

$$\begin{aligned}
\forall o, o' \in N : \quad o \rightarrow o' &\stackrel{\text{def}}{\iff} (\exists f \in \text{Fld} : o \xrightarrow{f} o' \in E) \vee (\exists d \in \text{Dir} : o \xrightarrow{d} o' \in C) \\
\forall v \in V, o \in N : \quad v \rightarrow o &\stackrel{\text{def}}{\iff} v \rightarrow o \in R.
\end{aligned}$$

The transitive closure of the points-to relation is denoted by  $\rightarrow^+$ :

$$\begin{aligned}
\forall o, o' \in N : \quad o \rightarrow^+ o' &\stackrel{\text{def}}{\iff} \exists o_1, \dots, o_k \in N : o \rightarrow o_1 \wedge o_1 \rightarrow o_2 \wedge \dots \wedge o_k \rightarrow o' \\
\forall v \in V, o \in N : \quad v \rightarrow^+ o &\stackrel{\text{def}}{\iff} \exists o_1, \dots, o_k \in N : v \rightarrow o_1 \wedge o_1 \rightarrow o_2 \wedge \dots \wedge o_k \rightarrow o.
\end{aligned}$$

The following equation decides whether two variables are *may-aliases* (i.e., possible aliases); it is defined over  $V \times V$ :

$$\text{MayAlias}(v, w) \stackrel{\text{def}}{\iff} (v = w) \vee (\exists o \in N : v \rightarrow o \wedge w \rightarrow o).$$

In some cases, reference sets in object nodes may help computing *must-aliases*, since a variable  $v$  ought to refer some object abstracted by a node  $n_X^H$  as soon as  $v \in X$ :

$$\text{MustAlias}(v, w) \stackrel{\text{def}}{\iff} (\forall n_X^H \in N : \{v, w\} \subset X \vee \{v, w\} \cap X = \emptyset).$$

*Reachability* is the reflexive and transitive closure of the points-to relation, it is defined over  $(V \cup N) \times N$ :

$$\text{Reach}(o, o') \stackrel{\text{def}}{=} o = o' \vee o \rightarrow^+ o'.$$

*Connectivity* is defined over  $(V \cup N) \times (V \cup N)$ . We take care of not considering paths which traverse incompatible nodes:

$$\begin{aligned} \text{Connect}(o_1, o'_1) &\stackrel{\text{def}}{\iff} o_1 = o'_1 \\ &\vee \left( \exists o_2, \dots, o_k, o'_2, \dots, o'_{k'} \in N : o_k = o_{k'} \wedge o_1 \rightarrow o_2 \wedge \dots \wedge o_{k-1} \rightarrow o_k \wedge o'_1 \rightarrow o'_2 \wedge \dots \wedge o'_{k'-1} \rightarrow o_{k'} \right. \\ &\quad \left. \wedge (\forall 1 \leq i < k, 1 \leq i' < k' : \text{Compat}(o_i, o'_{i'})) \right). \end{aligned}$$

Combness of individual connections is directly captured within points-to edges, but deciding combness of a full container is more technical (see Figure 1 for a simple example) and several competing definitions exist. A container object (resp. summary) node  $l$  abstracts a *comb* container (resp. containers) if two distinct nodes attached to  $l$  may not be connected, if each node attached to  $l$  is pointed by a comb edge, and if all outgoing edges of nodes *reachable* from  $l$  are also comb:

$$\begin{aligned} \text{FullComb}(l) &\stackrel{\text{def}}{\iff} \neg(\exists o, o' \in N, d \in \text{Dir} : \text{Reach}(l, o) \wedge o \xrightarrow{d} o' \in C \wedge \neg \text{Comb}(d)) \\ &\quad \wedge \neg(\exists o, o' \in N, d, d' \in \text{Dir} : l \xrightarrow{d} o \in C \wedge l \xrightarrow{d'} o' \in C \wedge \text{Connect}(o, o') \wedge \neg \text{Coincide}(d, d')). \end{aligned}$$

This definition of combness is very strong, because it relies on predicate *Connect*. A weaker definition is also useful, it is called *shallow combness* and considers first level pointers only:

$$\text{ShallowComb}(l) \stackrel{\text{def}}{\iff} \neg(\exists o \in N, d \in \text{Dir} : l \xrightarrow{d} o \in C \wedge \neg \text{Comb}(d)).$$

Combness and/or shallow combness are critical for many optimizations and (dependence) analyses. Some applications have been discussed in another paper [22] and many other are left for future work.

## 8 Putting the Analysis to Work

Many important issues have not been covered by the previous formalization and algorithm.

### 8.1 Cost Improvements

A naming scheme based on sets may induce an exponential number of nodes, but this happens only when many reference assignments are guarded with conditional expressions [18], and when many different iterator variables are attached to the same container. Edges labels can also hold an exponential number of primitive directions, especially when “flattening” directions. Still, it seems that container traversals associated with common algorithms exhibit much regularity, and we believe the size of direction expressions is less likely to grow large than the number of graph nodes.

Nevertheless, some approximations may be necessary to reduce the analysis worst case complexity: bounding the number of nodes via  $k$ -limiting on variable sets [18] is one of these. One may also consider approximate simplification rules for directions, either based on  $k$ -limiting or specific properties of direction operators. Such heuristics are easy to implement and to prove.

Dealing with standard containers provides a novel and more evolved way to deal with complexity issues. Roughly speaking, we propose a programmer-driven heuristic to switch between Chase-Wegman-Zadeck’s efficiency [1]—fast but lacking precision on destructive updates—and Sagiv-Reps-Wilhelm’s expensive materialization capabilities—with potentially exponential complexity. Growing variable sets on variable assignments is indeed critical for destructive updating of recursive structures: it enables killing on assignments of

the form  $v.f = w$  and on `delete()` or `put()` calls. But assuming that typed containers are used, one may choose to precisely kill *container* edges only. Then, variables would be added to reference and history sets for the `get()` method only. This heuristic dramatically reduces points-to graph growth while preserving a good precision. Compared to traditional cost-reduction heuristics which consider all pointers equally important, higher-level objects guide the compiler into more customizable behavior. Our technique unveils practical benefits of high-level (object-oriented) descriptions for static analysis and optimization.

## 8.2 Precision Improvements

One important feature of [18] has been left out of this presentation because it was of little use in the context of typed container structures, assuming that recursive pointer dereference will be much less frequent. The *sharing* property captures whether an object abstracted by a summary node may be targeted by two separate pointers or not. It is critical for precise materialization from summary nodes with “looping edges” and for shape classification (e.g., acyclic lists, trees, acyclic graphs, cyclic lists, and cyclic graphs). Sagiv and al. define a node predicate named *is\_shared* to capture sharing properties [18]: if  $is\_shared(n) = false$  then a single concrete object abstracted by  $n$  at a given run-time instance may only be pointed through one incoming points-to edge. Predicate *is\_shared* has also been extended to avoid confusing multiple fields [4]. Improving our analysis along these lines raises no fundamental problem but makes transformers difficult to understand and does not seem necessary for programs with explicit containers.

Finally, since we used no allocation point in the node naming scheme, an intermediate representation in Static Single Assignment (SSA) [6] would be preferred for the best precision.

## 9 Related Work

There are several methods addressing the shape and/or alias analysis problems. Some of them are based on program *annotations*, as in Hendren, Hummel and Nicolau [12]. Other methods approximate *access paths* induced by pointers to capture may-alias pairs, like Landi, Ryder and Zhang [16] or Deutsch [8]; Hendren and Nicolau use *path matrices* [13] to classify data structure shapes.

Analyses closer to our work use points-to graphs to capture more general kinds of data structure properties. Some of these are *store-based*—they use allocation sites to name graph—including Chase, Wegman and Zadeck [1] or Jones and Muchnick [14]. The latter uses multiple graphs per program point to cover all possible pointer configurations. Other schemes use *k-limiting* to bound the number of nodes, including Larus and Hilfinger [17], but usually suffer from exponential growth of the graph size.

Our work is based on Sagiv, Reps and Wilhelm *store-less shape graphs* and *variable set* node naming [18]. It still has a worst case exponential number of nodes but this only happens in tricky conditional pointer assignments. This naming covers multiple pointer configurations in a single graph and capture must-points-to properties for stack variables, critical for precise handling of destructive updating. We extended the naming with *history sets* to improve separation of incompatible pointer configurations in *summary* nodes. Sagiv, Reps and Wilhelm also introduced the *node materialization* scheme to “extract” objects from summary nodes. We generalized it to materialization from *object* nodes (for `get()` operations) allowing simultaneous coverage of multiple container configurations. We showed in Section 8.1 that dealing with explicit container types allowed further reduction of analysis complexity, compared to equivalent pointer-chasing implementations.

Corbera, Asenjo and Zapata [4] proposed additional improvements to shape graphs for the purpose of dependence analysis. Sagiv, Reps and Wilhelm even proposed a general *parametric* shape analysis framework [20] to enable systematic and safe extensions with so-called *instrumentation predicates*. But some of our extensions *cannot* be expressed using instrumentation predicates, e.g., materialization at `get()`. Most of these improvements are compatible and independent from our node naming and container-specific features; they could be useful for further precision and/or cost improvements.

Finally, our experiments established that converting container code to pointer-chasing syntax may dramatically decrease precision; extensions proposed in [18, 4, 20] are of little help, especially in the case of multiple simultaneous container traversals. Other precise techniques like Deutsch’s symbolic access paths [8] (based on Parikh mappings and regular expression decompositions) would also fail for lack of absolute naming of iterator positions.

## 10 Conclusion and Future Work

As programmers look forward to apply object-oriented design models in high-performance applications, compilers have to rely on higher level analyses and optimizations. Pointer analysis for container-centric applications is one of these, and we showed the benefit of using standard libraries and toolkits when exploiting abstract semantics of container and iterator methods. As a tradeoff between efficiency and precision, we proposed natural extensions to previous successful works, introducing combness and dedicated points-to graphs with iterator information.

Theoretical and experimental results confirm the design choices and applicability of our technique. Much work is left for further studies, including complete characterization and proof in abstract interpretation, context-sensitive and modular interprocedural extensions, and evolution of the prototype into a larger compiler infrastructure for real-sized experiments.

It would also be profitable to consider iterator interactions in a more evolved framework: current analysis is rather conservative in the case of iterator cloning. Comparing container positions could involve expressive abstractions such as Presburger arithmetics, especially for arrays [22]. Such approaches seem promising ways to further improve precision; the drawbacks being complexity and serious difficulties to handle structural updates like dynamic insertions and deletions.

**Acknowledgments** This work is supported by a CNRS-UIUC exchange program. We would like to thank Paul Feautrier for initiating research on combness and exchanging ideas.

## References

- [1] D. Chase, M. Wegman, and F. Zadeck. Analysis of pointers and structures. In *ACM Symp. on Principles of Programming Languages (PoPL'93)*, pages 296–310, Charleston, South Carolina, USA, Jan. 1993.
- [2] B.-C. Cheng and W.-M. Hwu. Modular interprocedural pointer analysis using access paths: Design, implementation, and evaluation. In *ACM Symp. on Programming Language Design and Implementation (PLDI'00)*, June 2000.
- [3] J.-D. Choi, M. Burke, and P. Carini. Efficient flow-sensitive interprocedural computation of pointer-induced aliases and side effects. In *20<sup>th</sup> ACM Symp. on Principles of Programming Languages (PoPL'93)*, pages 232–245, Charleston, South Carolina, USA, Jan. 1993.
- [4] F. Corbera, R. Asenjo, and E. Zapata. New shape analysis techniques for automatic parallelization of c codes. In *ACM Int. Conf. on Supercomputing (ICS'99)*, Rhodes, Greece, June 1999.
- [5] P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction of approximation of fixpoints. In *4<sup>th</sup> ACM Symp. on Principles of Programming Languages*, pages 238–252, Los Angeles, California, USA, Jan. 1977.
- [6] R. Cytron, J. Ferrante, B. K. Rosen, M. N. Wegman, and F. K. Zadeck. Efficiently computing static single assignment form and the control dependence graph. *ACM Trans. on Programming Languages and Systems*, 13(4):451–490, Oct. 1991.
- [7] B. A. Davey and H. A. Priestley. *Introduction to Lattices and Order*. Cambridge University Press, 1990.
- [8] A. Deutsch. Interprocedural may-alias analysis for pointers: beyond  $k$ -limiting. In *ACM Symp. on Programming Language Design and Implementation (PLDI'94)*, pages 230–241, Orlando, Florida, USA, June 1994.
- [9] J. Dolby and A. A. Chien. An evaluation of automatic object inline allocation techniques. In *ACM Symp. on Programming Language Design and Implementation (PLDI'00)*, Vancouver, British Columbia, Canada, June 2000.

- [10] R. Ghiya and L. J. Hendren. Connection analysis: A practical interprocedural heap analysis for C. *Int. Journal of Parallel Programming*, 24(6), 1996.
- [11] R. Ghiya and L. J. Hendren. Is it a tree, a dag, or a cyclic graph? A shape analysis for heap-directed pointers in C. In *23<sup>rd</sup> ACM Symp. on Principles of Programming Languages (PoPL'96)*, pages 1–15, St. Petersburg Beach, Florida, USA, Jan. 1996.
- [12] L. J. Hendren, J. Hummel, , and A. Nicolau. Abstractions for recursive pointer data structures: improving the analysis and transformation of imperative programs. In *ACM Symp. on Programming Language Design and Implementation (PLDI'92)*, pages 249–260, San Francisco, California, USA, June 1992.
- [13] L. J. Hendren and A. Nicolau. Parallelizing programs with recursive data structures. *IEEE Trans. on Parallel and Distributed Systems*, 1(1):35–47, Jan. 1990.
- [14] N. D. Jones and S. S. Muchnick. A flexible approach to interprocedural data flow analysis and programs with recursive data structures. *ACM Press*, 1982.
- [15] J. B. Kam and J. D. Ullman. Monotone data flow analysis frameworks. *Acta Informatica*, 7:309–317, 1977.
- [16] W. A. Landi, B. G. Ryder, and S. Zhang. Interprocedural modification side effect analysis with pointer aliasing. In *ACM Symp. on Programming Language Design and Implementation (PLDI'93)*, pages 56–67, Albuquerque, New Mexico, USA, June 1993.
- [17] J. R. Larus and P. N. Hilfinger. Detecting conflicts between structure accesses. In *ACM Symp. on Programming Language Design and Implementation (PLDI'88)*, pages 21–34, 1988.
- [18] M. Sagiv, T. Reps, and R. Wilhelm. Solving shape-analysis problems in languages with destructive updating. *ACM Trans. on Programming Languages and Systems*, 20(1):1–50, Jan. 1998.
- [19] S. Sagiv, T. W. Reps, and R. Wilhelm. Solving shape-analysis problems in languages with destructive updating. In *23<sup>rd</sup> ACM Symp. on Principles of Programming Languages (PoPL'96)*, pages 16–31, St. Petersburg Beach, Florida, USA, Jan. 1996.
- [20] S. Sagiv, T. W. Reps, and R. Wilhelm. Parametric shape analysis via 3-valued logic. In *26<sup>th</sup> ACM Symp. on Principles of Programming Languages (PoPL'99)*, pages 105–118, San Antonio, Texas, USA, Jan. 1999.
- [21] R. Wilson and M. S. Lam. Efficient context-sensitive pointer analysis for c programs. In *ACM Symp. on Programming Language Design and Implementation (PLDI'95)*, pages 1–12, La Jolla, California, USA, June 1995.
- [22] P. Wu, P. Feautrier, D. Padua, and Z. Sura. Points-to analysis for java with applications to loop optimizations. Technical Report CSRD TR-1585, U. of Illinois, Nov. 2000.
- [23] P. Wu and D. Padua. Containers, on the parallelization of general-purpose java programs. *Int. Journal of Parallel Programming*, 28(6):589–605, dec 2000.

## Appendix

### A Working With Directions

One may easily check the following orderings on direction expressions.

$$d'' \sqsubseteq d \parallel d' \sqsubseteq d \ast d' \quad \emptyset \sqsubseteq d \sqsubseteq d'' \sqsubseteq \ast''$$

In addition,  $d \sqsubseteq \ast$  if and only if  $\gamma(d)$  holds only *injective* mappings.

Here are the most important simplifications on direction expressions, according to equivalence relation  $\equiv$  (proofs of non-trivial results are provided in Section B.2).

- For a given iterator  $i$ , primitive direction  $\diamond_i$  only abstracts mappings whose domain is a singleton. This constraint induces strong simplifying properties:

$$\diamond_i^\# \equiv \diamond_i \# \diamond_i \equiv \diamond_i \parallel \diamond_i \equiv \diamond_i \quad (3)$$

$$\forall d, d' \in \text{Dir} : (\diamond_i \cdot d) \# (\diamond_i \cdot d') = (\diamond_i \cdot d) \parallel (\diamond_i \cdot d'). \quad (4)$$

- Primitive directions of an iterator  $i$  partition the set of *comb* connections: they are *exclusive* with respect to operator  $\cdot$  and composing them with  $\parallel$  yields all comb connections (or all connections when applying unary operator  $\#$ ).

$$\begin{aligned} \diamond_i \cdot \triangleleft_i &\equiv \diamond_i \cdot \triangleleft_i^\# \equiv \diamond_i \cdot \triangleright_i \equiv \diamond_i \cdot \triangleright_i^\# \equiv \triangleleft_i \cdot \triangleright_i \equiv \triangleleft_i^\# \cdot \triangleright_i \equiv \triangleleft_i \cdot \triangleright_i^\# \equiv \triangleleft_i^\# \cdot \triangleright_i^\# \equiv \emptyset \\ \triangleleft_i \parallel \diamond_i \parallel \triangleright_i &\equiv * \quad (\triangleleft_i \parallel \diamond_i \parallel \triangleright_i)^\# \equiv *^\# \end{aligned}$$

As a result,  $d \cdot (\triangleleft_i \parallel \triangleright_i)^\#$  abstract every connection in  $\gamma(d)$  whose domain does not contain the current position of  $i$ :

$$\gamma(d \cdot (\triangleleft_i \parallel \triangleright_i)^\#) = \{(\mu, \nu) \in \gamma(d) \mid \nu(i) \notin \text{Domain}(\mu)\}.$$

- Direction  $\emptyset$  is the zero of  $\cdot$  and the neutral element of  $\parallel$  and  $\#$ :  $a \parallel \emptyset \equiv a \# \emptyset \equiv a$  and  $a \cdot \emptyset \equiv \emptyset$ . Conversely, direction  $*^\#$  is the neutral element of  $\cdot$  and the zero of  $\parallel$  and  $\#$ :  $a \cdot *^\# \equiv a$  and  $a \parallel *^\# \equiv a \# *^\# \equiv *^\#$ .
- Unary operator  $\#$  is idempotent:

$$\forall d \in \text{Dir} : (d^\#)^\# \equiv d^\#.$$

It is *not* the case of binary operator  $\#$ .

- Remember  $\cdot$  and  $\parallel$  are the meet and join operators in  $\text{Dir}$ . Therefore, considering  $b \sqsubseteq a$ ,  $a \cdot b \equiv b$  and  $a \parallel b \equiv a$ . In particular,  $\cdot$  and  $\parallel$  are idempotent.
- A weaker property holds for binary operator  $\#$ . It states that unary operator  $\#$  is equivalent to iterated application of its binary counterpart.

$$\forall d \in \text{Dir} : \bigcup_{n \geq 1} \underbrace{d \# d \# \dots \# d}_n \equiv d^\#. \quad (5)$$

- Occurrences of  $\parallel$  can be replaced by  $\#$  under unary operator  $\#$ :

$$\forall a, b \in \text{Dir} : (a \parallel b)^\# \equiv (a \# b)^\#. \quad (6)$$

- Combness propagates through operator  $\cdot$ :

$$\forall a, b, c \in \text{Dir} : \text{Comb}(a) \implies a \cdot (b \# c) \equiv a \cdot (b \parallel c). \quad (7)$$

As a result,  $d$  is a comb direction if and only if it can be expressed without operator  $\#$  (neither binary nor unary versions). An inductive application of (7) gives a *linear-time algorithm* to check any direction expression for combness.

The four operators on directions have distributive and partial distributive properties (proofs in Section B.3).

$$\begin{aligned} a \cdot (b \parallel c) &\equiv (a \cdot b) \parallel (a \cdot c) & a \parallel (b \cdot c) &\equiv (a \parallel b) \cdot (a \parallel c) \\ a \cdot (b \# c) &\equiv (a \cdot b) \# (a \cdot c) & a \# (b \cdot c) &\equiv (a \# b) \cdot (a \# c) \\ a \parallel (b \# c) &\sqsubseteq (a \parallel b) \# (a \parallel c) & a \# (b \parallel c) &\equiv (a \# b) \parallel (a \# c) \\ (a \cdot b)^\# &\equiv a^\# \cdot b^\# & (a \# b)^\# &\equiv a^\# \# b^\# \end{aligned}$$

Distributivity properties of  $\cdot$  and  $\parallel$  (resp.  $\#$ ) are similar to those for logical operators: the lattice of directions is actually distributive. Moreover, (6) states that  $(a \parallel b)^\#$  can be simplified into  $(a \# b)^\#$ . For implementation and theoretical purposes, this may lead to the definition of *normal forms* for directions. We will only use the weaker “flattening” result, however.

**Theorem 3** Any direction  $d$  can be written as an equivalent expression where  $\cdot$  only operates on primitive directions  $\diamond_i$ ,  $\triangleleft_i$ ,  $\triangleright_i$ , and their unary  $\varkappa$  counterparts  $\triangleleft_i^\varkappa$  and  $\triangleright_i^\varkappa$  (for all iterators  $i$ ), where  $\varkappa$  operates on such  $\cdot$ -factors, and where  $\parallel$  operates on such  $\varkappa$ -terms.

Let us end this section with distributivity properties for  $\alpha$  and  $\gamma$ . By definition,  $\gamma$  is distributive over the meet operator of each lattice. Monotonicity only yields a partial result for  $\alpha$ :

$$\begin{aligned} \forall d, d' \in \text{Dir} : \gamma(d \cdot d') &= \gamma(d) \cap \gamma(d') \\ \forall C, C' \subseteq \text{Conn} : \alpha(C \cap C') &\sqsubseteq \alpha(C) \cdot \alpha(C'). \end{aligned}$$

Conversely,  $\alpha$  is distributive over the join operator of each lattice (this is proven in Section B.5). Monotonicity only yields a partial result for  $\gamma$ :

$$\begin{aligned} \forall C, C' \subseteq \text{Conn} : \alpha(C \cup C') &\equiv \alpha(C) \parallel \alpha(C') \\ \forall d, d' \in \text{Dir} : \gamma(d \parallel d') &\supseteq \gamma(d) \cup \gamma(d') \end{aligned}$$

## B Abstract Interpretation Proofs

Two simple properties of  $\gamma$  are used in the proofs.

**Lemma 2** In the definitions of  $\gamma(e \parallel e')$  and  $\gamma(e \varkappa e')$ , one may choose  $\mu_e$  and  $\mu'_e$  with disjoint domains without lack of generality.

**Lemma 3** If  $(\mu, \nu)$  belongs to  $\gamma(e)$  and  $\mu_{/P}$  is the restriction of  $\mu$  to a set  $P$  of positions (possibly empty), then  $(\mu_{/P}, \nu)$  also belongs to  $\gamma(e)$ .

### B.1 Combness and Comparison

We first review combness properties. All primitive directions abstract only connections  $(\mu, \nu)$  where  $\mu$  is injective;  $\gamma(d \cdot d')$  holds only injective mappings  $\mu$  if and only if  $\gamma(d)$  and/or  $\gamma(d')$  do. The disjoint range condition in  $\gamma(d \parallel d')$  ensures that it holds injective mappings  $\mu$  when both  $\gamma(d)$  and  $\gamma(d')$  do, and the condition is also necessary.

Abstract interpretation allows us to give a direct definition of Coincide and Disjoint. It is easy to check that the following equations are equivalent to the inductive definition in Section 3.1.

$$\begin{aligned} \text{Coincide}(d, d') &\iff \forall (\mu, \nu) \in \gamma(d), (\mu', \nu) \in \gamma(d') : \text{Domain}(\mu) \subset \text{Domain}(\mu') \vee \text{Domain}(\mu') \subset \text{Domain}(\mu) \\ \text{Disjoint}(d, d') &\iff \forall (\mu, \nu) \in \gamma(d), (\mu', \nu) \in \gamma(d') : \text{Domain}(\mu) \cap \text{Domain}(\mu') = \emptyset. \end{aligned}$$

Then, it is easy to see that  $\text{Coincide}(d, d')$  is equivalent to  $d \parallel d' \sqsubseteq \diamond_i$  for some iterator  $i$ , and  $\text{Disjoint}(d, d')$  is equivalent to  $d \cdot d' \equiv \emptyset$ .

Remember that coincident directions abstract connections whose domain is either empty or a singleton. This is implied by Lemma 3:  $\mu$  and  $\mu'$  with empty or singleton domains are the only kind of mapping which satisfy the definition of Coincide.

### B.2 Simplifications

- Equivalence of  $*$  and  $\triangleleft_i \parallel \diamond_i \parallel \triangleright_i$ .

$$\begin{aligned} \gamma(\triangleleft_i \parallel \diamond_i \parallel \triangleright_i) &= \{(\mu, \nu) \in \text{Conn} \mid \exists \mu_1, \mu_2, \mu_3 : \text{Injective}(\mu_1) \wedge \text{Injective}(\mu_3) \\ &\quad \wedge (\forall p \in \text{Domain}(\mu) : \mu(p) \in \{\mu_1(p), \mu_2(p), \mu_3(p)\}) \\ &\quad \wedge (\forall p \in \text{Domain}(\mu_1) : p < \nu(\mathbf{i})) \\ &\quad \wedge \text{Domain}(\mu_2) = \{\nu(\mathbf{i})\} \\ &\quad \wedge (\forall p \in \text{Domain}(\mu_3) : p > \nu(\mathbf{i})) \\ &\quad \wedge \text{Range}(\mu_1) \cap \text{Range}(\mu_2) = \emptyset \\ &\quad \wedge \text{Range}(\mu_2) \cap \text{Range}(\mu_3) = \emptyset \\ &\quad \wedge \text{Range}(\mu_1) \cap \text{Range}(\mu_3) = \emptyset\}. \end{aligned}$$

Thus, the three mappings  $\mu_1, \mu_2, \mu_3$  are required to have disjoint domains, to be injective, and to have disjoint ranges. This is equivalent to  $\mu$  being injective itself:

$$\gamma(\langle \cdot \rangle_i \parallel \diamond_i \parallel \triangleright_i) = \{(\mu, \nu) \in \text{Conn} : \text{Injective}(\mu)\}.$$

- Iterated application of operator  $\varkappa$ . Let  $(\mu, \nu) \in \gamma(d^n)$ .

$$\exists \mu' : \text{Domain}(\mu') = \text{Domain}(\mu) \wedge (\mu', \nu) \in \gamma(d).$$

Therefore, applying Lemma 3,

$$\forall p \in \text{Domain}(\mu), \exists \mu_p : \mu(p) = \mu_p(p) \wedge (\mu_p, \nu) \in \gamma(d).$$

Calling  $n$  the cardinal of  $\text{Domain}(\mu)$ ,

$$\exists (\mu_k)_{1 \leq k \leq n} : (\forall p \in \text{Domain}(\mu), \exists 1 \leq k \leq n : \mu(p) = \mu_k(p)) \wedge (\forall 1 \leq k \leq n : (\mu_k, \nu) \in \gamma(d)),$$

hence  $(\mu, \nu) \in \gamma\left(\bigcup_{n \geq 1} \underbrace{d \varkappa d \varkappa \cdots \varkappa d}_n\right)$ . The reciprocal is straightforward.

### B.3 Distributivity and Partial Distributivity

Let  $a, b$  and  $c$  be any directions. We first prove that  $a \cdot (b \varkappa c) \equiv (a \cdot b) \varkappa (a \cdot c)$ .

$$\begin{aligned} (\mu, \nu) \in \gamma(a \cdot (b \varkappa c)) &\iff \exists \mu_b, \mu_c : (\forall p \in \text{Domain}(\mu) : \mu(p) \in \{\mu_b(p), \mu_c(p)\}) \\ &\quad \wedge (\mu, \nu) \in \gamma(a) \wedge (\mu_b, \nu) \in \gamma(b) \wedge (\mu_c, \nu) \in \gamma(c) \\ &\iff (\exists \mu_b, \mu_c : (\forall p \in \text{Domain}(\mu) : \mu(p) \in \{\mu_b(p), \mu_c(p)\}) \\ &\quad \wedge (\mu_b, \nu) \in \gamma(a) \cap \gamma(b) \wedge (\mu_c, \nu) \in \gamma(a) \cap \gamma(c)) \\ &\iff (\mu, \nu) \in \gamma((a \cdot b) \varkappa (a \cdot c)). \end{aligned}$$

We prove that  $a \varkappa (b \cdot c) \equiv (a \varkappa b) \cdot (a \varkappa c)$ . Let  $(\mu, \nu) \in \gamma(a \varkappa (b \cdot c))$ . There are mappings  $\mu_a$  and  $\mu_{bc}$  such that

$$(\forall p \in \text{Domain}(\mu) : \mu(p) \in \{\mu_a(p), \mu_{bc}(p)\}) \wedge (\mu_a, \nu) \in \gamma(a) \wedge (\mu_{bc} \in \gamma(b) \cap \gamma(c)).$$

Mappings  $\mu_a$  and  $\mu_{bc}$  satisfy both conditions of  $a \varkappa b$  and  $a \varkappa c$ , thus  $(\mu, \nu) \in \gamma((a \varkappa b) \cdot (a \varkappa c))$ . Conversely, if  $(\mu, \nu) \in \gamma((a \varkappa b) \cdot (a \varkappa c))$ , there are mappings  $\mu_{ab}, \mu_{ac}, \mu_b$ , and  $\mu_c$  such that

$$\begin{aligned} (\forall p \in \text{Domain}(\mu) : \mu(p) \in \{\mu_{ab}(p), \mu_b(p)\}) \wedge \mu(p) \in \{\mu_{ac}(p), \mu_c(p)\}) \\ \wedge (\mu_a, \nu) \in \gamma(a) \wedge (\mu'_a, \nu) \in \gamma(a') \wedge (\mu_b \in \gamma(b)) \wedge (\mu_c \in \gamma(c)). \end{aligned}$$

We can build  $\mu_a$  such that  $\mu_a(p) = \mu_{ab}(p)$  when  $\mu(p) = \mu_{ab}(p)$  and  $\mu_a(p) = \mu_{ac}(p)$  when  $\mu(p) = \mu_{ac}(p)$ , and we then replace both  $\mu_{ab}$  and  $\mu_{ac}$  in the previous result:

$$\begin{aligned} (\forall p \in \text{Domain}(\mu) : \mu(p) \in \{\mu_a(p), \mu_b(p)\}) \wedge \mu(p) \in \{\mu_a(p), \mu_c(p)\}) \\ \wedge (\mu_a, \nu) \in \gamma(a) \wedge (\mu_b \in \gamma(b)) \wedge (\mu_c \in \gamma(c)). \end{aligned}$$

Now, when  $\mu(p)$  differs from  $\mu_a(p)$ , it has to be equal to both  $\mu_b(p)$  and  $\mu_c(p)$ . We may thus build a mapping  $\mu_{bc}$  such that

$$(\forall p \in \text{Domain}(\mu) : \mu(p) \in \{\mu_a(p), \mu_{bc}(p)\}) \wedge (\mu_a, \nu) \in \gamma(a) \wedge (\mu_{bc} \in \gamma(b) \cap \gamma(c)),$$

hence  $(\mu, \nu) \in \gamma(a \varkappa (b \cdot c))$ .

Proofs for mutual distributivity between  $\cdot$  and  $\parallel$  are likewise: the disjoint ranges condition does not interfere with the former reasoning.



We prove that  $a \parallel (b \text{ ; } c) \sqsubseteq (a \parallel b) \text{ ; } (a \parallel c)$ . Let  $(\mu, \nu) \in \gamma(a \parallel (b \text{ ; } c))$ . There are mappings  $\mu_a$  and  $\mu_b$  and  $\mu_c$  such that

$$\begin{aligned} (\forall p \in \text{Domain}(\mu) : \mu(p) \in \{\mu_a(p), \mu_b(p), \mu_c(p)\}) \wedge (\mu_a, \nu) \in \gamma(a) \wedge (\mu_b, \nu) \in \gamma(b) \wedge (\mu_c, \nu) \in \gamma(c) \\ \wedge \text{Range}(\mu_a) \cap (\text{Range}(\mu_b) \cup \text{Range}(\mu_c)) = \emptyset. \end{aligned}$$

Since  $\text{Range}(\mu_a) \cap (\text{Range}(\mu_b) \cup \text{Range}(\mu_c)) = (\text{Range}(\mu_a) \cap \text{Range}(\mu_b)) \cup (\text{Range}(\mu_a) \cap \text{Range}(\mu_c))$ , it shows that  $(\mu, \nu) \in \gamma((a \parallel b) \text{ ; } (a \parallel c))$ .

We eventually prove that  $a \text{ ; } (b \parallel c) \equiv (a \text{ ; } b) \parallel (a \text{ ; } c)$ . Let  $(\mu, \nu) \in \gamma(a \text{ ; } (b \parallel c))$ . There are mappings  $\mu_{ab}$ ,  $\mu_{ac}$ ,  $\mu_b$  and  $\mu_c$  such that

$$\begin{aligned} (\forall p \in \text{Domain}(\mu) : \mu(p) \in \{\mu_{ab}(p), \mu_{ac}(p), \mu_b(p), \mu_c(p)\}) \wedge (\mu_{ab}, \nu) \in \gamma(a) \wedge (\mu_{ac}, \nu) \in \gamma(a) \\ \wedge (\mu_b, \nu) \in \gamma(b) \wedge (\mu_c, \nu) \in \gamma(c) \wedge (\text{Range}(\mu_{ab}) \cup \text{Range}(\mu_b)) \cap (\text{Range}(\mu_{ac}) \cup \text{Range}(\mu_c)) = \emptyset. \end{aligned}$$

We can build  $\mu_a$  such that  $\mu_a(p) = \mu_{ab}(p)$  when  $\mu(p) = \mu_{ab}(p)$  and  $\mu_a(p) = \mu_{ac}(p)$  when  $\mu(p) = \mu_{ac}(p)$ , and we then replace both  $\mu_{ab}$  and  $\mu_{ac}$  in the previous result:

$$\begin{aligned} (\forall p \in \text{Domain}(\mu) : \mu(p) \in \{\mu_a(p), \mu_b(p), \mu_c(p)\}) \wedge (\mu_a, \nu) \in \gamma(a) \wedge (\mu_b, \nu) \in \gamma(b) \wedge (\mu_c, \nu) \in \gamma(c) \\ \wedge (\text{Range}(\mu_a) \cup \text{Range}(\mu_b)) \cap (\text{Range}(\mu_a) \cup \text{Range}(\mu_c)) = \emptyset. \end{aligned}$$

It shows that  $(\mu, \nu) \in \gamma(a \text{ ; } (b \parallel c))$ . Conversely, let  $(\mu, \nu) \in \gamma(a \text{ ; } (b \parallel c))$ . There are mappings  $\mu_a$  and  $\mu_b$  and  $\mu_c$  such that

$$\begin{aligned} (\forall p \in \text{Domain}(\mu) : \mu(p) \in \{\mu_a(p), \mu_b(p), \mu_c(p)\}) \wedge (\mu_a, \nu) \in \gamma(a) \wedge (\mu_b, \nu) \in \gamma(b) \wedge (\mu_c, \nu) \in \gamma(c) \\ \wedge \text{Range}(\mu_b) \cap \text{Range}(\mu_c) = \emptyset. \end{aligned}$$

Let  $\mu_{ab}$  (resp.  $\mu_{ac}$ ) be the restriction of  $\mu_a$  to positions  $p$  where  $\mu_a(p) \notin \text{Range}(\mu_c)$  (resp.  $\mu_a(p) \notin \text{Range}(\mu_b)$ ). Since  $\text{Range}(\mu_b) \cap \text{Range}(\mu_c) = (\text{Range}(\mu_{ab}) \cup \text{Range}(\mu_b)) \cap (\text{Range}(\mu_{ac}) \cup \text{Range}(\mu_c))$ ,  $(\mu, \nu) \in \gamma(a \text{ ; } (b \parallel c))$ .

## B.4 Lattice Completeness

First of all, we prove completeness for the meet operator. Let  $e$  be a *flat* expression associated with some direction  $d$  (see Theorem 3). If the *unary* mix operator is not used in  $e$ , any *descending* chain starting from  $e$  is finite: it traverses sub-expressions (without unary mix operators) only. Now, consider an expression  $e$  with unary mix operators. If there is an infinite descending chain starting from  $e$ , (5) combined with the previous result enforces that there is an expression  $e'$  in that chain with less occurrences of the unary mix operator. An induction on the number of occurrences of the unary mix operator yields a contradiction. This proves that every descending chain is finite. Therefore,  $\text{Dir}$  is a complete meet semi-lattice.

Eventually, Theorem 2.16 in [7] shows that  $\text{Dir}$  is a complete lattice, because it also has a top element ( $*$ ).

## B.5 Galois Connection

Let us prove (2). Consider a set  $C \in \text{Conn}$ , and call  $d = \alpha(C)$ . For all  $c \in C$ , definition of  $\alpha$  enforces that  $\alpha(c) \sqsubseteq d$ , and then  $c \in \gamma(d)$  by monotonicity of  $\gamma$ .

Conversely, considering a direction  $d$ , there is a direction  $d'$  such that  $\alpha(\gamma(d)) = \sqsubseteq d$ . Monotonicity of  $\gamma$  yields  $\gamma(\alpha(\gamma(d))) \subseteq \gamma(d)$ . But considering  $C = \gamma(d)$  in the result of the previous paragraph yields  $\gamma(d) \subseteq \gamma(\alpha(\gamma(d)))$ . Therefore  $\gamma(d) = \gamma(\alpha(\gamma(d)))$ . Injectivity of  $\gamma$  terminates the proof.

Eventually, let us prove that  $\alpha$  is distributive over the join operator of each lattice. Consider two sets of connections  $C$  and  $C'$ . Monotonicity ensures that  $\alpha(C \cup C') \sqsubseteq \alpha(C) \parallel \alpha(C')$ . Conversely, let  $(\mu, \nu) \in \gamma(\alpha(C) \parallel \alpha(C'))$ , and consider  $\mu_1, \mu_2$  such that  $\forall p \in \text{Domain}(\mu) : \mu(p) \in \{\mu_1(p), \mu_2(p)\}$ ,  $(\mu_1, \nu) \in \gamma(\alpha(C))$ ,  $(\mu_2, \nu) \in \gamma(\alpha(C'))$ , and  $\text{Range}(\mu_1) \cap \text{Range}(\mu_2) = \emptyset$ . Using monotonicity again, both  $(\mu_1, \nu)$  and  $(\mu_2, \nu)$  belong to  $\gamma(\alpha(C \cup C'))$ . Moreover,  $\mu_1$  and  $\mu_2$  can be chosen with disjoint domains (from Lemma 2). Since  $\mu_1$  and  $\mu_2$  also have disjoint ranges, a simple inductive reasoning on direction expressions shows that  $(\mu, \nu)$  also belongs to  $\gamma(\alpha(C \cup C'))$ . Therefore,  $\gamma(\alpha(C) \parallel \alpha(C')) \subseteq \gamma(\alpha(C \cup C'))$ , which concludes the proof.

## B.6 Local Safety

We prove that transfer functions of iterator and container related operations are safe with respect to container connections. This proof relies on the hypothesis that other parts of the analysis—not related to container edges—are correct.

Let  $(\mu, \nu)$  be a connection between a container abstracted by a node  $l$  and objects abstracted by a node  $o$ . Consider a container edge  $l \xrightarrow{d} o$  in a points-to graph  $G$ , such that  $(\mu, \nu) \in \gamma(d)$ . When  $\text{Domain}(\mu) = \emptyset$ , there might be no such edge, but we do as if there was one labeled with  $d = \emptyset$ . Consider an operation whose concrete and abstract transfer functions are  $\tau$  and  $\tau^\sharp$ , respectively. We show that  $\alpha(\tau(\mu, \nu)) \sqsubseteq \tau^\sharp(\alpha(\mu, \nu))$  for methods `first()`, `last()`, `advance()`, `retreat()`, `delete()`, `insert()`, `put()`, and `get()`.

**i. `l.first()` (`l.last()`).** We first consider a container `l` which is not abstracted by node  $l$  (the source of the container edge). Concrete transfer function  $\tau$  is defined as  $\tau(\mu, \nu) = (\mu, \nu')$  where  $\nu'(\mathbf{i})$  is undefined and  $\nu'(j) = \nu(j)$  for all  $j \neq \mathbf{i}$ . Moreover,  $\tau^\sharp$  replaces primitive directions of  $\mathbf{i}$  with  $*$ , hence  $\alpha(\mu, \nu') \sqsubseteq \tau^\sharp(\alpha(\mu, \nu')) \equiv \tau^\sharp(\alpha(\mu, \nu))$ .

If `l` is abstracted by node  $l$ ,  $\tau(\mu, \nu) = (\mu, \nu'')$  where  $\nu'(\mathbf{i}) = 0$  (instead of being undefined) and  $\nu'(j) = \nu(j)$  for all  $j \neq \mathbf{i}$ . Because the former definition of  $\nu'$  is the restriction of  $\nu''$  to a smaller set of iterators, one has  $\alpha(\mu, \nu'') \sqsubseteq \alpha(\mu, \nu') \sqsubseteq \tau^\sharp(\alpha(\mu, \nu))$ .

**i. `advance()` (`retreat()`).** Concrete semantics of `advance()` states that  $\tau(\mu, \nu) = (\mu, \nu')$  where  $\nu'(\mathbf{i}) = \nu(\mathbf{i}) + 1$  and  $\nu'(j) = \nu(j)$  for all  $j \neq \mathbf{i}$ . Therefore  $(\mu, \nu') \in \gamma(\alpha(\mu, \nu)[\langle \mathbf{i} / \diamond_{\mathbf{i}}, * / \triangleright_{\mathbf{i}} \rangle])$ , hence  $\alpha(\tau(\mu, \nu)) \sqsubseteq \tau^\sharp(\alpha(\mu, \nu))$ .

**i. `delete()`.** From the concrete semantics of `delete()`,  $\tau(\mu, \nu) = (\mu', \nu)$  where  $\mu'(\nu(\mathbf{i}))$  is undefined and  $\mu'(p) = \mu(p)$  for all  $p \neq \nu(\mathbf{i})$ . As a result,  $(\mu', \nu) \in \gamma(\langle \mathbf{i}^* \triangleright_{\mathbf{i}} \rangle)$ , hence  $\alpha(\tau(\mu, \nu)) \sqsubseteq \alpha(\mu, \nu) \cdot \langle \mathbf{i}^* \triangleright_{\mathbf{i}} \rangle \equiv \tau^\sharp(\alpha(\mu, \nu))$ .

**i. `insert(v)`.** We suppose that variable  $v$  references the  $n$ -th object created during program execution. From the concrete semantics of `insert()`,  $\tau(\mu, \nu) = (\mu', \nu)$  where  $\mu'(\nu(\mathbf{i}) + 1) = n$ ,  $\mu'(p) = \mu(p)$  for all  $p \leq \nu(\mathbf{i})$ , and  $\mu'(p + 1) = \mu(p)$  for all  $p > \nu(\mathbf{i})$ . Of course,  $(\mu, \nu)$  belongs to  $\gamma(\alpha(\mu, \nu) \# \triangleright_{\mathbf{i}})$  as well as  $(\mu'', \nu)$  where  $\mu''$  is the restriction of  $\mu'$  to  $\{\nu(\mathbf{i}) + 1\}$ . Therefore,  $\alpha(\tau(\mu, \nu)) \sqsubseteq \alpha(\mu, \nu) \# \triangleright_{\mathbf{i}} \equiv \tau^\sharp(\alpha(\mu, \nu))$ .

**i. `put(v)`.** Consider the first rule in the transfer function, where the container edge points to a node  $o'$  which does not abstract the object referenced by  $v$ . Then, the concrete semantics of `put()` coincides with this of `delete()`.

The next two rules are very similar to the `insert()` method. We suppose again that variable  $v$  references the  $n$ -th object created during program execution:  $\tau(\mu, \nu) = (\mu', \nu)$  where  $\mu'(\nu(\mathbf{i})) = n$  and  $\mu'(p) = \mu(p)$  for all  $p \neq \nu(\mathbf{i})$ . Then,  $(\mu, \nu)$  belongs to  $\gamma(\alpha(\mu, \nu) \# \diamond_{\mathbf{i}})$  as well as  $(\mu'', \nu)$ , where  $\mu''$  is the restriction of  $\mu'$  to  $\{\nu(\mathbf{i})\}$ . Notice  $\alpha(\mu, \nu) \equiv \emptyset$  in the second rule and  $\alpha(\mu, \nu) = d$  in the third one. Thus, in both cases,  $\alpha(\mu, \nu) \# \diamond_{\mathbf{i}} \equiv \tau^\sharp(\alpha(\mu, \nu))$ .

Remains the fourth rule, where  $o$ , the edge's target, abstracts the object referenced by  $v$ , and where  $\diamond_{\mathbf{i}} \in d$  (i.e.,  $\diamond_{\mathbf{i}}$  appears in every expression of  $d$ ). The transfer function is  $\tau^\sharp(d) = d \# (d \cdot \diamond_{\mathbf{i}})$ . Once again,  $\tau(\mu, \nu) = (\mu', \nu)$  where  $\mu'(\nu(\mathbf{i})) = n$  is the creation number associated with the single object abstracted by object node  $o$ , and  $\mu'(p) = \mu(p)$  for all  $p \neq \nu(\mathbf{i})$ . From Lemma 1 (proven at the end of this section), there is a connection  $(\mu_{\mathbf{i}}, \nu) \in \gamma(d)$  associated with a real execution of the program which satisfies  $\text{Domain}(\mu_{\mathbf{i}}) = \{\nu(\mathbf{i})\}$ . Of course,  $\mu_{\mathbf{i}}$  can be choosen such that  $\mu_{\mathbf{i}}(\nu(\mathbf{i})) = n$ . This shows that  $\mu'$  is partitionned into a mapping  $\mu_{\mathbf{i}}$ , and a mapping  $\mu''$ , defined as the restriction of  $\mu$  to positions  $p \neq \nu(\mathbf{i})$ . These two mappings have disjoint domains but their ranges may actually intersect. If they do not intersect, it means that  $\mu''$  has empty domain, because there is only one object abstracted by the target node; hence  $\mu' = \mu''$  and  $(\mu', \nu) \in \gamma(d)$ . On the contrary,  $(\mu'', \nu) \in \gamma(d)$  and  $(\mu_{\mathbf{i}}, \nu) \in \gamma(d \cdot \diamond_{\mathbf{i}})$ , thus  $(\mu', \nu) \in \gamma(d \# (d \cdot \diamond_{\mathbf{i}}))$ . Finally,  $\alpha(\tau(\mu, \nu)) = \alpha(\mu', \nu) \sqsubseteq d \# (d \cdot \diamond_{\mathbf{i}}) = \tau^\sharp(d)$ .

**v = i. `get()`.** The first rule (after strong nullification) does not interfere with container edges.

Consider the edge created after materialization in the second and third rules. These rules only differ in the replacement of  $d$  by  $d \cdot \diamond_{\mathbf{i}}$  when  $d$  is a comb direction. The materialized container edge targets

an object node  $o = n_{X \cup \{v\}}^{H \cup \{v\}}$ , and we call  $n$  the single object abstracted by  $o$ . From the concrete semantics of `get()`,  $\tau(\mu, \nu) = (\mu', \nu)$  where  $\mu'$  is the restriction of  $\mu$  to  $\{p : \mu(p) = \mu(\nu(i))\}$ . Of course, having  $\tau^\sharp(d) = d$  in the non-comb case is safe, because  $\text{Domain}(\mu') \subseteq \text{Domain}(\mu)$  and thus  $\alpha(\mu', \nu) \sqsubseteq \alpha(\mu, \nu)$ . On the other hand, suppose that  $d$  is a comb direction. This ensures injectivity of  $\mu$ , hence  $\mu(p) = \nu(i) \implies p = \nu(i)$ . Therefore, the domain of  $\mu'$  is the singleton  $\{\nu(i)\}$ . Finally, we have  $(\mu', \nu) \in \gamma(d)$  and  $(\mu', \nu) \in \gamma(\diamond_i)$ , hence  $\alpha(\tau(\mu, \nu)) \sqsubseteq \alpha(\mu, \nu) \cdot \diamond_i \equiv \tau^\sharp(\alpha(\mu, \nu))$ .

This section ends with the proof of Lemma 1. We follow the inductive definition of direction expressions, using operators  $\cdot$ ,  $\parallel$  and  $\#$ .

First of all, any connection  $(\mu, \nu)$  such that  $\text{Domain}(\mu) = \{\nu(i)\}$  is abstracted by  $\diamond_i \parallel d$  and  $\diamond_i \# d$ . Moreover,  $(\mu, \nu) \in \gamma(d) \implies (\mu, \nu) \in \gamma(d^\#)$ . Eventually, consider a direction of the form  $\diamond_i \cdot d$ . As noticed in Section 5.4.3, only two rules—in the transfer functions for `get()` and `put()`—build directions of the form  $\diamond_i \cdot d$ . Other graph transformations, including join at control-flow merge points, do not build such directions. Consider a container edge labeled  $\diamond_i \cdot d$  at a program point  $p$  following a call to `get()` or `put()`. Let  $(\mu', \nu) \in \gamma(\diamond_i \cdot d)$  be a connection associated with a real execution of the program, i.e., in the collecting semantics at  $p$ . From the construction of  $\mu'$  in the previous safety proofs,  $\nu(i)$  belongs to  $\text{Domain}(\mu')$ . Lemma 3 terminates the proof.

## B.7 Merging Edges

We now prove Theorem 2. Let  $G = (V, N, R, E, C)$  and  $G' = (V', R', N', E', C')$  be two points-to graphs associated with adjacent incoming control-flow edges. Supposing that  $G$  and  $G'$  are conservative approximations of points-to relations, it must be proven that  $G'' = (V'', R'', N'', E'', C'') = G \sqcup G'$  is still a conservative approximation for all paths leading to the merge point through the chosen incoming control-flow edges.

The case of variables, objects and their connections is straightforward, since  $V'' = V \cup V'$ ,  $N'' = N \cup N'$ ,  $R'' = R \cup R'$  and  $E'' = E \cup E'$ .

Let us prove that  $C'' = C \cup C'$  is a conservative approximation of container connections after the merge point. Considering in turn conditionals and loops, we show that  $e'' = l \xrightarrow{d \parallel d'} o$  abstracts every connection previously abstracted by either  $e = l \xrightarrow{d} o \in C$  or  $e' = l \xrightarrow{d'} o \in C'$ .

- Exclusion rules on conditionals (with or without `else` case, and arbitrary `switch` constructs) only allows a single execution path to be valid for a given instance of the merge point during execution. Choosing  $d \parallel d'$  is thus safe.
- Loops induce two merge points: we consider a given *run-time instance* of either the exit point or the iteration point. Let  $(\mu, \nu) \in \gamma(d)$  and  $(\mu', \nu) \in \gamma(d')$  be two configurations abstracted by  $e$  and  $e'$ , for this precise run-time instance. Monotonicity properties of the iterative analysis enforce that  $\mu \subseteq \mu'$  or  $\mu' \subseteq \mu$ , i.e., the domain of one mapping includes the other and both mappings coincide on the smaller domain. As a result, joining  $\mu$  and  $\mu'$  into a new mapping  $\mu''$  yields either  $\mu'' = \mu$  or  $\mu'' = \mu'$ . The mapping  $\mu''$  from the configuration  $(\mu'', \nu'') \in \gamma(e'')$  at the same run-time instance is thus partitioned into  $\mu' \setminus \mu$  and  $\mu \setminus \mu'$  with disjoint ranges (one of them is empty). Hence  $(\mu'', \nu'') \in \gamma(d \parallel d')$ .

This concludes the proof, conditionals and loops being the only cases of merge points (no `gotos`).



---

Unité de recherche INRIA Rocquencourt  
Domaine de Voluceau - Rocquencourt - BP 105 - 78153 Le Chesnay Cedex (France)  
Unité de recherche INRIA Lorraine : LORIA, Technopôle de Nancy-Brabois - Campus scientifique  
615, rue du Jardin Botanique - BP 101 - 54602 Villers-lès-Nancy Cedex (France)  
Unité de recherche INRIA Rennes : IRISA, Campus universitaire de Beaulieu - 35042 Rennes Cedex (France)  
Unité de recherche INRIA Rhône-Alpes : 655, avenue de l'Europe - 38330 Montbonnot-St-Martin (France)  
Unité de recherche INRIA Sophia Antipolis : 2004, route des Lucioles - BP 93 - 06902 Sophia Antipolis Cedex (France)

---

Éditeur  
INRIA - Domaine de Voluceau - Rocquencourt, BP 105 - 78153 Le Chesnay Cedex (France)  
<http://www.inria.fr>  
ISSN 0249-6399