

Bottom-up Context-Sensitive Pointer Analysis for Java

Yu Feng Xinyu Wang Isil Dillig Thomas Dillig

University of Texas at Austin
{ yufeng, xwang, isil, tdillig }@cs.utexas.edu

ABSTRACT

This paper describes a new bottom-up, subset-based, and context-sensitive pointer analysis for Java. The main novelty of our technique is the constraint-based handling of virtual method calls and instantiation of method summaries. Since our approach generates polymorphic method summaries, it can be context-sensitive without reanalyzing the same method multiple times. We have implemented this algorithm in a tool called SCUBA, and we compare it with k -CFA and k -obj algorithms on Java applications from the DaCapo and Ashes benchmarks. Our results show that the new algorithm achieves better or comparable precision to k -CFA and k -obj analyses at only a fraction of the cost.

1. INTRODUCTION

Pointer analysis is a key enabling technology underlying many program analysis, software engineering, and compiler optimization tasks. Given a pointer variable p , pointer analysis statically determines the set of all heap objects that p may point to. The result of such an analysis can be used to resolve important program analysis questions, such as whether two pointers can be aliases or whether a heap location may be referenced in a given piece of code.

While existing pointer analysis algorithms differ along many dimensions, a key feature that determines the precision of an algorithm is *context-sensitivity*. In particular, a context-sensitive analysis respects the call/return semantics of procedure calls and does not yield spurious points-to facts that arise from interprocedurally unrealizable paths. Furthermore, a context-sensitive analysis distinguishes between heap objects that are allocated at the same program location, but due to different invocations of the same method. While more precise than context-insensitive ones, context-sensitive algorithms are much harder to scale to real programs, and many existing techniques use approximations of full context-sensitivity. For instance, *object-sensitive* analyses [13, 14, 23] only distinguish callsites where the receiver objects are different, and k -CFA analyses [22, 11] differentiate contexts by tracking callstrings up to some fixed length k .

Context-sensitivity can be achieved either by performing a *top-down* or *bottom-up* interprocedural analysis. Top-down analyses start at entry methods of a program and analyze callers before callees. In contrast, bottom-up analyses start at leaf methods of the callgraph and analyze callees before callers. Since top-down algorithms analyze every method in a known calling context, they are simpler to design and im-

plement, but they need to re-analyze the same method several times for different calling contexts. In contrast, bottom-up interprocedural analyses generate a *polymorphic* method summary that may be used in *any* calling context to get context-sensitive results. While generating a polymorphic points-to summary is trickier than determining points-to information at a particular call site, bottom-up analyses do not need to reanalyze the same method several times¹ and have the potential to scale better. In addition, the results of a bottom-up pointer analysis are *reusable*: For instance, using a bottom-up pointer analysis, we can analyze a library just once and reuse its summary for many different clients.

In this paper, we present a bottom-up context- and field-sensitive pointer analysis algorithm for Java. A key novel feature of our approach is the constraint-based treatment of virtual method calls. Similar to many other approaches, our method starts with an imprecise callgraph and refines the callgraph as points-to facts are discovered. However, we construct the callgraph in a purely bottom-up fashion by predicating points-to facts on the possible dynamic types of the receiver object. As method summaries are propagated up the call chain, these dynamic types are resolved, thereby allowing the refutation of infeasible call targets and spurious points-to facts in a context-sensitive manner.

Another salient feature of our approach is that it can generate polymorphic method summaries without performing expensive case splits on possible aliasing patterns at call sites. In particular, a key challenge in bottom-up pointer analysis is how to generate method summaries that soundly capture the aggregate effect of a call to method m under *any* possible aliasing relation at m 's call sites. Most previous techniques deal with this difficulty either by performing case splits on all possible aliasing patterns [6, 3] (which can cause exponential blow-up) or by using unification-based methods [9, 28, 12, 27] (which are imprecise compared to subset-based methods). A main advantage of our technique is that it is as precise as subset-based methods despite modeling the unknown state of the heap in a simple and uniform way. In particular, since our technique does not perform *strong updates*² to heap locations, it can soundly account for the callee's side effects by performing a fixed-point computation during summary instantiation.

¹In bottom-up algorithms, the reanalysis of a method is necessary only when it is part of an SCC in the callgraph.

²A strong update to memory location o kills the existing points-to facts for o , while a *weak update* does not.

<pre> class X { Z f; Z g; void bar(Z z) { this.f = z; } } </pre>	<pre> class A { X x; X y; void a1() { y = new Y(); x = y; Z z = new Z(); foo(z); } void a2() { x = new X(); y = x; Z z = new Z(); foo(z); } void foo(Z a) { x.bar(a); y.bar(a); } } </pre>
<pre> class Y extends X { void bar (Z z) { this.g = z; } } </pre>	

Figure 1: Code example to illustrate our approach

We have implemented our algorithm in a tool called SCUBA³, and we compare its scalability and precision with top-down pointer analysis algorithms implemented in CHORD [15]. Our experimental results on programs from the DaCapo and Ashes benchmark suites indicate that SCUBA achieves better or comparable precision to k -CFA and k -object-sensitive algorithm at a fraction of the cost.

To summarize, this paper makes the following contributions:

- We present a bottom-up, subset-based, and context-sensitive pointer analysis for Java. A key novelty of our approach is the handling of virtual method calls using constraint-based techniques.
- We describe a new method for summarizing and instantiating points-to facts. Unlike previous techniques, our approach does not case-split on aliasing patterns and guarantees soundness by performing fixed-point computation during summary instantiation.
- We describe an implementation of our algorithm and compare it with k -CFA and k -obj algorithms on benchmarks taken from the DaCapo and Ashes test suites.

2. EXAMPLE

This section illustrates our analysis using an example that showcases virtual method calls and the need for context-sensitivity. Consider the code shown in Figure 1, which defines classes X , Y , and A . Here, Y is a subclass of X and overrides X 's `bar` method. Class A has two instance variables x and y of type X . The `a1` and `a2` methods of A initialize the x and y fields and call method `foo`, which invokes `bar` on both x and y . For concreteness, suppose we want to know whether $x.f$ and $y.g$ can be aliases at the end of `a1` and `a2`.

Our algorithm starts by analyzing $X::\text{bar}$ and $Y::\text{bar}$ since these methods are leaf procedures in the callgraph. The summaries for $X::\text{bar}$ and $Y::\text{bar}$ are shown in Figure 2 (a) and (b). We depict both method summaries as well as local points-to facts in the form of a graph, where nodes correspond to abstract heap objects and directed edges denote may-point-to relations. Our method summaries only include points-to edges that may be *added* due to an invocation of the summarized method. In particular, method summaries do not include points-to relations that already exist on method entry, and since our analysis does not apply strong updates, no points-to edges can be removed as

³SCUBA stands for SCAlable UPwards Bottom-up Analysis.

a result of analyzing a method. Hence, the summary for a method m can be thought of as a bag of (symbolic) points-to edges that are introduced due to an invocation of m .

Now, consider the summary for $X::\text{bar}$ shown in Figure 2(a). We use arg_0 to denote the `this` pointer, and arg_1 , arg_2 to denote the first and second formal parameters. For parameter i , we use $\text{arg}_i.\epsilon$ to denote the heap object pointed to by arg_i on method entry. Hence, according to the summary in Figure 2(a), $X::\text{bar}$ adds a points-to edge via field f from the heap object pointed to by the `this` pointer to the heap object pointed to by the first parameter of method $X::\text{bar}$. Note that at a call site of $X::\text{bar}$, parameter z may have multiple points-to targets, hence, the single edge in $X::\text{bar}$'s summary may actually introduce multiple points-to edges at a call site of $X::\text{bar}$. The summary for $Y::\text{bar}$, which is shown in Figure 2(b), is very similar to that of $X::\text{bar}$. The only difference is that the points-to edge is from the g field of $\text{arg}_0.\epsilon$ to $\text{arg}_1.\epsilon$ since $Y::\text{bar}$ stores into the g field.

Now consider method $A::\text{foo}$, whose summary is shown in Figure 2(c). Since x has type X , the call `x.bar(a)` could either invoke $X::\text{bar}$ or $Y::\text{bar}$. Hence, to analyze the method call, we instantiate the summaries of both $X::\text{bar}$ and $Y::\text{bar}$, but the points-to edges induced by instantiating method $T::\text{bar}$ are qualified by a constraint that stipulates that the dynamic type of the receiver is T .

As an example, consider the potential target $X::\text{bar}$ of the call `x.bar(a)`. Here, the only points-to edge in the summary for $X::\text{bar}$ is from the f field of $\text{arg}_0.\epsilon$ to $\text{arg}_1.\epsilon$. Since arg_0 of $X::\text{bar}$ corresponds to `this.x` at the call site, $\text{arg}_0.\epsilon$ instantiates to $\text{arg}_0.\epsilon.x$, which denotes the memory location pointed to by `this.x`. On the other hand, arg_1 in $X::\text{bar}$ corresponds to parameter z at the call site, hence $\text{arg}_1.\epsilon$ translates to a location with the same name, i.e., $\text{arg}_1.\epsilon$. Therefore, as shown in Figure 2(c), instantiating the summary of $X::\text{bar}$ for this call site induces an edge from $\text{arg}_0.\epsilon.x$ to $\text{arg}_1.\epsilon$, but this edge is qualified by the constraint $\text{type}(\text{arg}_0.\epsilon.x) = X$. Similarly, instantiating the summary of $Y::\text{bar}$ for the call `x.bar(a)` induces a points-to edge from the g field $\text{arg}_0.\epsilon.x$ to $\text{arg}_1.\epsilon$, but qualified by the constraint $\text{type}(\text{arg}_0.\epsilon.x) = Y$.

Continuing with the method invocation `y.bar(a)` in the second line of $A::\text{foo}$, we again instantiate the summaries of $X::\text{bar}$ and $Y::\text{bar}$, since the type of y is X . However, this time, the location named $\text{arg}_1.\epsilon$ used in the summaries of

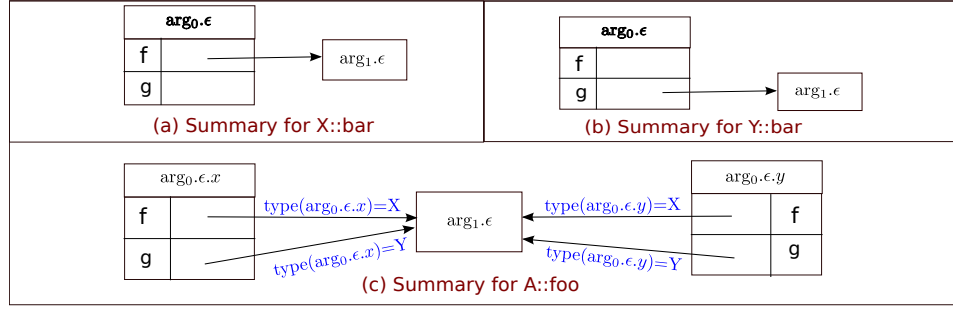


Figure 2: Method summaries computed by our algorithm

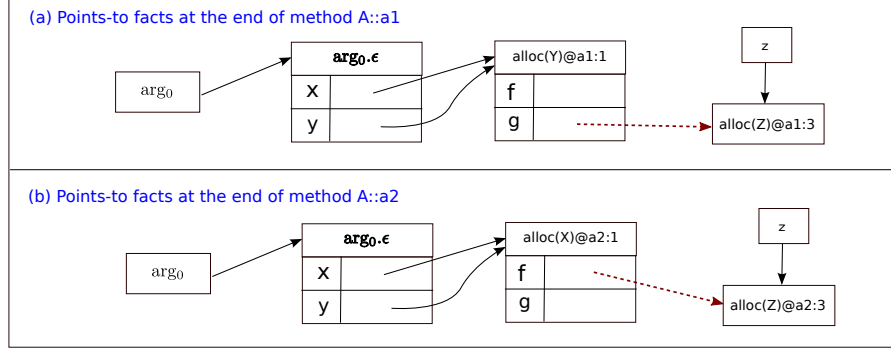


Figure 3: Points-to facts at the end of a1, a2. Dashed lines indicate points-to edges added after the foo call.

$X::\text{bar}$ and $Y::\text{bar}$ correspond to the location pointed to by `this.y`, which is denoted by $\text{arg}_0.\epsilon.y$ in method `foo`. Hence, as shown in Figure 2(c), the summary for `foo` includes points-to edges from the `f` and `g` fields of $\text{arg}_0.\epsilon.y$ to $\text{arg}_1.\epsilon$, again qualified by the appropriate type constraints. Observe that our approach is context-sensitive because the two different invocations of `bar` induce different points-to relations. Also, even though `this.x` and `this.y` may alias at a call site of `foo`, we represent their points-to targets on method entry using two *separate* locations called $\text{arg}_0.\epsilon.x$ and $\text{arg}_0.\epsilon.y$. As we will see shortly, this approach is sound as long as we do not perform strong updates because, in contexts where `this.x` and `this.y` alias, $\text{arg}_0.\epsilon.x$ and $\text{arg}_0.\epsilon.y$ will instantiate to the same location.

Now consider method `a1` defined in `A`. The *solid* black edges in Figure 3(a) denote the points-to facts that hold right before the call to `foo`. In particular the location named $\text{arg}_0.\epsilon$ denotes the object pointed to by the `this` pointer, and both the `x` and `y` fields of $\text{arg}_0.\epsilon$ point to a location called $\text{alloc}(Y)@a1:1$, which corresponds to the memory allocated at the first line of method `a1`. Throughout the paper, we use the notation $\text{alloc}(T)@Ctx$ to denote heap objects of type `T` that are allocated in context `Ctx`, where `Ctx` is a list of method name and line number pairs. Similarly, in Figure 3(a), the local variable `z` points to a freshly allocated memory block denoted by $\text{alloc}(Z)@a1:3$.

We now turn to the method invocation `foo(z)` in `a1`. Here, the location named $\text{arg}_1.\epsilon$ in `foo`'s summary corresponds to the location pointed to by `z` in `a1`, which is $\text{alloc}(Z)@a1:3$. On the other hand, the locations named $\text{arg}_0.\epsilon.x$ and $\text{arg}_0.\epsilon.y$ in `foo` represent the locations pointed to by `this.x` and `this.y` in `a1` respectively. Following the chain of points-to edges in Figure 3(a), we see that $\text{arg}_0.\epsilon.x$ and $\text{arg}_0.\epsilon.y$ both

correspond to the location $\text{alloc}(Y)@a1:1$. Since this allocation is tagged with type `Y`, the constraints $\text{type}(\text{arg}_0.\epsilon.x) = X$ and $\text{type}(\text{arg}_0.\epsilon.y) = X$ both evaluate to false, while $\text{type}(\text{arg}_0.\epsilon.x) = Y$ and $\text{type}(\text{arg}_0.\epsilon.y) = Y$ evaluate to true. Hence, instantiating `foo`'s summary induces a single points-to edge from the `g` field of $\text{alloc}(Y)@a1:1$ to $\text{alloc}(Z)@a1:3$, which is shown with the dotted edge in Figure 3(a).

Now, when we instantiate `foo`'s summary in `a2`, $\text{arg}_0.\epsilon.x$ and $\text{arg}_0.\epsilon.y$ in the summary both correspond to $\text{alloc}(X)@a2:1$, as shown in Figure 3(b). Hence, constraints $\text{type}(\text{arg}_0.\epsilon.x) = X$ and $\text{type}(\text{arg}_0.\epsilon.y) = X$ evaluate to true, while constraints $\text{type}(\text{arg}_0.\epsilon.x) = Y$, $\text{type}(\text{arg}_0.\epsilon.y) = Y$ evaluate to false. Thus, instantiating `foo`'s summary in `a2` introduces a points-to edge from the `f` field of $\text{alloc}(X)@a2:1$ to $\text{alloc}(Z)@a2:3$, as shown with the dotted edge in Figure 3(b).

As we can see from Figure 3, the analysis determines that `x.f` and `y.g` are not aliases in either `a1` or `a2`. Observe that an analysis that is either context-insensitive or based on an imprecise callgraph would conclude otherwise. Also, even though there are two different calls to `foo` and four different calls to `bar`, our algorithm analyzes each method once.

3. CONCEPTUAL FOUNDATIONS

Before describing our analysis in detail, we first describe a conceptual framework that lays the foundations of our algorithm. The exposition in this section is largely conceptual in nature and does not describe an implementable analysis.

3.1 Preliminaries

We describe the main ideas underlying our technique using a may-points-to graph, which we refer to as an *abstract heap*:

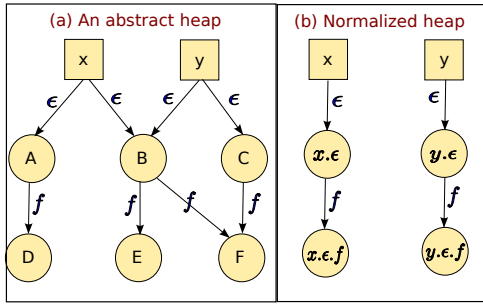


Figure 4: An abstract heap and its normal form

DEFINITION 1. (Abstract heap) An abstract heap H is a graph (N, E) where N is a set of nodes corresponding to abstract memory locations, and E is a set of directed edges between nodes labeled with field names or ϵ . An edge (o_1, o_2, f) indicates that the f field of o_1 may point to o_2 .

Here, an abstract memory location represents either the stack location of a variable or a set of heap objects. The edge label ϵ is used to model points-to relations from stack locations to heap objects. The root nodes of an abstract heap denote locations of variables, and we write $\text{root}(H)$ to indicate the set of root nodes of H . Given two abstract heaps H_1 and H_2 , $H_1 \cup H_2$ represents the abstract heap containing nodes and edges from both H_1 and H_2 . Given a heap H and edges E , we write $H \setminus E$ to denote the heap that contains all nodes and edges in H except the set of edges E .

3.2 Normalization of Abstract Heaps

Given an abstract heap H , we now define a *normalization operation* $\mathbb{N}(H)$, which yields a *normalized heap* H^* and a mapping ζ from nodes in H to sets of nodes in H^* .

DEFINITION 2. (Normal form) Given heap $H = (N, E)$, $\mathbb{N}(H)$ yields normalized heap $H^* = (N^*, E^*)$ and mapping $\zeta : N \rightarrow 2^{N^*}$ such that:

1. If $x \in \text{root}(H)$, then $x \in N^*$ and $\zeta(x) = \{x\}$.
2. If $(o, o', f) \in E$ and $o^* \in \zeta(o)$, then $o^*.f \in N^*$, $o^*.f \in \zeta(o')$, and $(o^*, o^*.f, f) \in E^*$.

We use the notation $\mathbb{N}(H) = H^*$ to indicate that H^* is in normal form, and we write $\text{Map}(H, H^*) = \zeta$ to indicate that ζ maps nodes of H to sets of nodes in H^* .

EXAMPLE 1. Consider the abstract heap H shown in Figure 4(a) and its normal form H^* in Figure 4(b). Here, $\text{Map}(H, H^*)$ yields the following mapping ζ :

$$\begin{aligned} \zeta(x) &= \{x\} & \zeta(y) &= \{y\} & \zeta(A) &= \{x.\epsilon\} \\ \zeta(B) &= \{x.\epsilon, y.\epsilon\} & \zeta(C) &= \{y.\epsilon\} & \zeta(D) &= \{x.\epsilon.f\} \\ \zeta(E) &= \{x.\epsilon.f, y.\epsilon.f\} & \zeta(F) &= \{x.\epsilon.f, y.\epsilon.f\} \end{aligned}$$

We use heap normal forms to model the heap on entry to a method. In particular, H^* corresponds to a “generic” heap representing the unknown points-to targets of object o ’s f field as $o.f$. While the mapping ζ from abstract heap H to its normal form H^* differs for each call site, the normalized heap for a method m is the same irrespective of calling context. Observe that no pair of abstract memory locations alias each other in a normalized heap, and every location has exactly one points-to target for a given field.

SUMANALYZE(H, S, A):

input: abstract heap H , code S , intraprocedural analysis A
output: abstract heap H'

```

(1) let  $H^* = \text{NormalForm}(H)$ 
(2) let  $H'^* = \text{Analyze}(H^*, S, A)$ 
(3) let  $\Delta = H'^* \setminus \text{default}(H'^*)$ 
(4) let  $\zeta_0 = \text{Map}(H, H^*)$ 
(5) let  $H' = H$ ; let  $\zeta = \zeta_0$ 
(6) do {
(7)    $\zeta_0 = \zeta$ 
(8)    $H' = \zeta^{-1}(\Delta) \cup H'$ 
(9)    $\zeta = \text{Map}(H', H^*)$ 
(10) }
(11) while( $\zeta \neq \zeta_0$ )
(12) return  $H'$ 

```

Figure 5: Basic structure of summary-based analysis

Given an abstract heap $H = (N, E)$ and its normal form $H^* = (N^*, E^*)$ such that $\text{Map}(H, H^*) = \zeta$, we define ζ^{-1} to be a mapping from N^* to 2^N such that $n \in \zeta^{-1}(n^*)$ if and only if $n^* \in \zeta(n)$. Observe that there is a many-to-many relationship between the nodes of an abstract heap and those of its corresponding normal form.

EXAMPLE 2. Consider the heap from Ex 1. We have:

$$\begin{aligned} \zeta^{-1}(x) &= \{x\} & \zeta^{-1}(y) &= \{y\} \\ \zeta^{-1}(x.\epsilon) &= \{A, B\} & \zeta^{-1}(y.\epsilon) &= \{B, C\} \\ \zeta^{-1}(x.\epsilon.f) &= \{D, E, F\} & \zeta^{-1}(y.\epsilon.f) &= \{E, F\} \end{aligned}$$

The mapping ζ^{-1} is important in summary-based analysis because it allows us to instantiate a method summary to a particular abstract heap at a call site. We use the notation ζ^{-1} to denote the extension of ζ^{-1} that maps any element that is not in the domain of ζ^{-1} to itself. Given heap $H^* = (N^*, E^*)$, we also write $\zeta^{-1}(H^*)$ to denote a heap $H = (N, E)$ where $(o_1, o_2, f) \in E$ iff there exists an edge $(o_1^*, o_2^*, f) \in E^*$ such that $o_1 \in \zeta^{-1}(o_1^*)$ and $o_2 \in \zeta^{-1}(o_2^*)$.

DEFINITION 3. (Default edge) We say an edge (n, n', f) is a default edge of an abstract heap if $n' = n.f$.

Given a heap H , we write $\text{default}(H)$ to denote the set of default edges in H .

3.3 Summary-Based Pointer Analysis

We now explain the basic idea underlying our summary-based analysis, assuming a family of pointer analyses that are sound and weakly-updating. Given code snippet S , an abstract heap H , and pointer analysis A , we write $H' = \text{Analyze}(H, S, A)$ to indicate that, if statement S is executed in an environment that satisfies abstract heap H , then analyzing code S using pointer analysis A yields a heap H' which conservatively models the concrete heap after S .

The basic structure of our summary-based pointer analysis is shown in Figure 5. The algorithm **SUMANALYZE** takes as input an abstract heap H , a code snippet S , and a weakly-updating pointer analysis A , and works as follows:

- Line (1) constructs the normalized heap H^* representing the unknown state of the heap before executing S
- Line (2) analyzes S without making any assumptions about points-to facts that hold before S
- Line (3) generates a polymorphic points-to summary Δ which characterizes side effects of code S

Program P $:= C^+$
 ClassDecl C $:= \text{class } T_1 [\text{extends } T_2]? \{F^*; M^*\}$
 FieldDecl F $:= T \text{ fld_name};$
 MethodDecl M $:= m(T_0 v_0, \dots, T_k v_k) = \{V^*; I;\}$
 VarDecl V $:= T \text{ var_name};$
 Instruction I $:= v_1 = v_2 \mid v_1 = v_2.f \mid v_1.f = v_2$
 $\mid v = \text{new}^\rho T \mid \text{if}(\ast) I_1 \text{ else } I_2 \mid I_1; I_2$
 $\mid m^\rho @ T(v_1, \dots, v_n) \mid v_0.m^\rho(v_1, \dots, v_n)$

Figure 6: Core language used for our formalization

- Lines (4)-(11) instantiate the summary Δ by performing a fixed-point computation
- H' at line (12) models the state of the heap after S when S is executed in an environment satisfying H

Before discussing details, let us first understand in what way this algorithm is “summary-based”. Since H^* can be constructed in a context-independent manner, we can analyze S in isolation and compute its side effects without knowing the points-to facts that hold before S . Hence, lines (1)-(3) in Figure 5 correspond to *summary generation*. On the other hand, lines (4)-(11) perform *summary instantiation* by computing the context-specific mapping ζ_0 and by adding all points-to edges that represent S ’s side effects.

The most involved part of the above algorithm is the fixed-point computation at lines (4)-(11). Intuitively, the algorithm maps each edge in the summary to a set of edges at the callsite by using the mapping ζ^{-1} and adds these edges to the initial abstract heap H . However, as new edges are added to H , the mapping ζ must to be recomputed since locations used in the summary may map to new additional locations after the summary has been applied. Hence, H' is recomputed until the mapping ζ from H' to H^* stabilizes.

It is easy to see that SUMANALYZE is sound because (i) the underlying pointer analysis does not apply strong updates to memory locations, and (ii) the summary is applied to a fixed-point. In particular, observe that H' overapproximates H and $\{H'\}S\{H'\}$ is a valid Hoare triple.

The reader may wonder why it is necessary to re-apply the summary until the mapping ζ reaches a fixed point in Figure 5. This is necessary because our summary Δ encodes all possible side effects of code snippet S , but not the *order* in which they happen. Hence, while the fixed point computation is required for soundness, an immediate corollary is that the procedure SUMANALYZE can be less precise than *Analyze* if the underlying pointer analysis A is flow-sensitive. This point is illustrated by the following example:

EXAMPLE 3. Consider the abstract heap shown in Figure 4(a) and the following code snippet S :

(1) $\mathbf{z} = \mathbf{y}.f$; (2) $\mathbf{x}.f = \text{new } \mathbf{G}$;

Here, a flow-sensitive analysis can determine that \mathbf{z} can only point to objects E and F from Figure 4(a), but not the \mathbf{G} object allocated at line (2). However, SUMANALYZE(H, S, A) would yield objects E, F, G as the points-to targets of \mathbf{z} .

4. FORMALIZATION OF ALGORITHM

While the previous section describes core ideas of the analysis, it omits many important details. This section describes our full algorithm using a core object-oriented language.

(Field selector) $\eta : f \mid \eta.f \mid \eta^*$
 (Heap obj) $o : a_i.\eta \mid \text{alloc}(T)@\bar{\rho}$
 (Abstract loc) $\pi : o \mid a_i \mid v_i@\bar{\rho}$
 (Pts set) $\theta : o \rightarrow \phi$
 (Abstract heap) $\Gamma : (\pi \times f) \rightarrow \theta$
 (Summaries) $\Upsilon : (T \times M) \rightarrow \Gamma$

Figure 7: Abstract domains used in our analysis

4.1 Language

We formalize our algorithm using the core language shown in Figure 6. Here, a program consists of one or more class declarations C , which introduces a class named T_1 with optional superclass T_2 . Instructions include assignments, loads, stores, heap allocations (marked with a unique program point ρ), non-deterministic conditionals, sequences, static method calls $m^\rho @ T(\dots)$ (also marked with program point ρ), and virtual method calls $v_0.m^\rho(\dots)$. We also assume that the first argument of a method is always the *this* pointer, and that, if a class T inherits a method m from its superclass T' , then there exists an implicit method declaration m in T with the same implementation as in T' .

4.2 Abstract Domains

Figure 7 shows the abstract domains that we need for describing our algorithm. Since our analysis is bottom-up, we differentiate between two kinds of heap objects o : *Access paths* of the form $a_i.\eta$ represent (caller-allocated) unknown heap objects reachable through the i ’th argument, whereas objects named $\text{alloc}(T)@\bar{\rho}$ represent heap objects of type T that are allocated either in the currently analyzed method or in a transitive callee. Specifically, $a_i.f_1\dots f_n$ denotes the unknown locations reachable on method entry through a series of field accesses $f_1\dots f_n$ from the i ’th argument. We use the notation $a_i.(f_1\dots f_n)^*$ to denote all unknown locations reachable from a_i through any combination of field selectors f_1, \dots, f_n . For instance, $a_0.(f.g)^*$ represents the infinite set of access paths $a_0.f, a_0.g, a_0.f.f, a_0.g.f$ and so on. As we will see in Section 4.4, access paths in our analysis correspond to node labels of the normalized heap from Section 3.

Abstract memory locations π are either heap objects o or stack locations. In particular, a_i denotes the stack location of the i ’th argument, and $v_i@\bar{\rho}$ denotes the location of local variable v_i under context $\bar{\rho}$. We represent calling contexts using a sequence of program points ρ_1, \dots, ρ_n , where each ρ_i corresponds to some call or allocation site. For instance, a memory location named $\text{alloc}(T)@\rho_1\rho_2$ corresponds to a heap object allocated at program point ρ_2 of some method m which is invoked at call site ρ_1 . Similarly, $v@\rho_1$ denotes the local variable v of some method m when m is invoked at callsite ρ_1 . Since our analysis builds contexts in a bottom-up way, local variables declared in the currently analyzed method m do not have any context information; hence, we abbreviate the locations of locals in the current method as v_i .

DEFINITION 4. (**Argument-derived location**) We say location π is derived from an argument, written $\text{arg}(\pi)$, if it is either (i) a_i representing the location of the i ’th argument or (ii) a heap object represented with an access path $a_i.\eta$.

Abstract heap Γ maps each field f of location π to a points-to set. A (guarded) points-to set θ is a set of pairs (o, ϕ) where o is a heap object and ϕ is a constraint. As discussed in Section 2, we use constraints to predicate points-to facts

on dynamic types of receivers. In particular, constraints ϕ are formulas in the theory of equality with uninterpreted functions, defined according to the following grammar:

Function f := pts | alloc | ς_i
 Term t := c | v | $f(\vec{t})$
 Formula ϕ := \top | \perp | $\text{type}(t) = T$ | $\phi_1 \wedge \phi_2$ | $\phi_1 \vee \phi_2$

Here, terms include constants c , variables v , and function applications $f(\vec{v})$ where f is either the binary function pts, alloc, or an n -ary function drawn from $\varsigma_1, \dots, \varsigma_k$. Formulas are composed of \top (true), \perp (false), and conjunctions and disjunctions of equality constraints of the form $\text{type}(t) = T$, where T is a type constant. In addition to the usual function and equality axioms, the alloc and type functions obey the additional axiom $\forall x. \text{type}(\text{alloc}(T, \vec{p})) = T$, which states that the type of an allocation of type T is T .

Since we will convert heap objects to terms, we define an operation called *lift*(π), abbreviated $\bar{\pi}$, as follows:

$$\begin{aligned} \bar{a}_i &= a_i & \overline{\text{alloc}(T)@p} &= \text{alloc}(T, \vec{p}) \\ \bar{\pi}.f &= \text{pts}(\bar{\pi}, f) & \bar{\pi}.(\vec{f})^* &= \varsigma_i(\bar{\pi}, \vec{f}) \text{ } (\varsigma_i \text{ fresh}) \end{aligned}$$

We also assume an operation $\text{lift}^{-1}(t)$ which is the inverse of *lift*(π). That is, given a term t , $\text{lift}^{-1}(t)$ yields an abstract memory location representation of that term.

EXAMPLE 4. The constraint $\text{type}(\overline{a.\epsilon.f}) = A \wedge \text{type}(\overline{a.\epsilon.f}) = B$ is unsatisfiable since the dynamic type of $a.\epsilon.f$ cannot simultaneously be A and B . On the other hand, the constraint $\text{type}(\overline{a.\epsilon.g^*}) = A \wedge \text{type}(\overline{a.\epsilon.g^*}) = B$, which translates to $\text{type}(\varsigma_1(\text{pts}(a, \epsilon), g)) = A \wedge \text{type}(\varsigma_2(\text{pts}(a, \epsilon), g)) = B$, is satisfiable because two distinct occurrences of $a.\epsilon.g^*$ may correspond to different objects (e.g., $a.\epsilon.g$ and $a.\epsilon.g.g$).

We now define a function *has_type*(θ, T) which generates a constraint that evaluates to true if *some* element in points-to set θ can have dynamic type T :

DEFINITION 5. *has_type*(θ, T) Given a points-to set θ , the function *has_type*(θ, T) yields the following constraint:

$$\bigvee_{(\pi_i, \phi_i) \in \theta} (((\text{type})(\bar{\pi}_i) = T) \wedge \phi_i)$$

As we will see in the next section, the *has_type* function will be useful for analyzing virtual method calls.

Now, going back to the abstract domains from Figure 7, the environment Υ maps each method M in class T to its corresponding summary, which is just an abstract heap Γ summarizing M 's side effects. Applying the summary at a call site allows us to determine points-to relations of the caller without having to reanalyze the callee. In addition, our method summaries include points-to information for locals in the summarized method. While the inclusion of local variables in the summary is not necessary for analyzing the caller, this design choice allows us to determine points-to sets for *all* program variables without employing a separate top-down pass. In particular, the summary of the program's entry method yields points-to sets for all program variables.

4.3 Operations on Abstract Domains

We now describe some operations on abstract domains that simplify the description of our algorithm. In particular, we

first need an operation that looks up points-to targets of locations from Γ . However, since our algorithm constructs the initial heap on method entry in a demand-driven way, the lookup operation retrieves so-called *default targets* for locations that are *not* yet in the domain of Γ :

DEFINITION 6. (Default target) Given an argument-derived location π and a field f , the default target of the f field of π , written $\text{def}(\pi, f)$, is given as follows:

$$\text{def}(\pi, f) = \begin{cases} \pi & \text{if } \pi = \pi'.(\vec{f})^* \text{ and } f \in \vec{f} & (1) \\ \pi'.(f.\vec{g})^* & \text{if } \pi = \pi'.f.\vec{g} & (2) \\ \pi.f & \text{otherwise} & (3) \end{cases}$$

In other words, if field f is *not* part of a recursive field cycle (line 3), then the default target for field f of an argument derived location π is $\pi.f$, just like the normal form heaps from Section 3. However, if f is part of a recursive field cycle, then our analysis collapses this cycle into a single abstract memory location (lines 1-2). For example, $\text{def}(a.\text{next}, \text{next}) = a.\text{next}^*$ (line 2), and $\text{def}(a.\text{next}^*, \text{next}) = a.\text{next}^*$ (line 1). The summarization of recursive field cycles into access paths of the form $a.f^*$ is needed to ensure termination of the fixed-point computation performed by our algorithm.

We now define a *field lookup* operation:

DEFINITION 7. (Field look-up) Given heap Γ , field f , and location π , the field lookup operation $\Gamma[\pi, f]$ retrieves the points-to target for π 's f field:

$$\Gamma[\pi, f] = \begin{cases} \Gamma(\pi, f) \cup \{(\text{def}(\pi, f), \top)\} & \text{if } \arg(\pi) \\ \Gamma(\pi, f) & \text{otherwise} \end{cases}$$

Since our algorithm does not explicitly add default edges to the abstract heap, $\Gamma[\pi, f]$ always yields $\text{def}(\pi, f)$ as part of the points-to set of $\pi.f$ if π is an argument derived location. Now, since our analysis performs weak updates, we need to merge two points-to sets using the following join operator:

DEFINITION 8. (Join of points-to sets) The join operator \sqcup on points-to sets is defined as follows:

$$(\theta_1 \sqcup \theta_2)(o) = \begin{cases} \theta_1(o) \vee \theta_2(o) & \text{if } o \in \text{dom}(\theta_1) \cap \text{dom}(\theta_2) \\ \theta_1(o) & \text{if } o \in \text{dom}(\theta_1) \text{ and } o \notin \text{dom}(\theta_2) \\ \theta_2(o) & \text{if } o \in \text{dom}(\theta_2) \text{ and } o \notin \text{dom}(\theta_1) \end{cases}$$

Observe that if an object o is in both points to sets θ_1 and θ_2 , we take the disjunction of the constraints associated with o . We also extend this join operator to abstract heaps in the expected way. That is, for a location π and field f , $(\Gamma_1 \sqcup \Gamma_2)(\pi, f)$ yields $\Gamma_1(\pi, f) \sqcup \Gamma_2(\pi, f)$.

In our analysis, we sometimes need to predicate points-to information on constraints. For this purpose, we define an operation $\theta \downarrow \phi$ which conjoins ϕ with every constraint in θ :

DEFINITION 9. (Projection of θ on ϕ)

$$\theta \downarrow \phi = \{(\pi_i, \phi_i \wedge \phi) \mid (\pi_i, \phi_i) \in \theta\}$$

Finally, we extend the field lookup operation on points-to sets as follows:

DEFINITION 10. (Field lookup for points-to set)

$$\Gamma[\theta, f] = \bigsqcup_{(\pi_i, \phi_i) \in \theta} \Gamma[\pi_i, f] \downarrow \phi_i$$

That is, $\Gamma[\theta, f]$ includes the points-to target of every element in θ under the appropriate constraints. This notation will become useful for describing the analysis of loads and stores.

$$\begin{aligned}
(1) \quad & \frac{\Gamma' = \Gamma[(v_1, \epsilon) \leftarrow \Gamma[v_2, \epsilon]]}{\Upsilon, \Gamma \vdash v_1 = v_2 : \Gamma'} \\
(2) \quad & \frac{\theta = \Gamma[v_2, \epsilon] \quad \Gamma' = \Gamma[(v_1, \epsilon) \leftarrow \Gamma[\theta, f]]}{\Upsilon, \Gamma \vdash v_1 = v_2.f : \Gamma'} \\
(3) \quad & \frac{\theta_1 = \Gamma[v_1, \epsilon] \quad \theta_2 = \Gamma[v_2, \epsilon] \quad \Gamma' = \Gamma[(o_i, f) \leftarrow (\Gamma(o_i, f) \sqcup (\theta_2 \downarrow \phi_i)) \mid (o_i, \phi_i) \in \theta_1]}{\Upsilon, \Gamma \vdash v_1.f = v_2 : \Gamma'} \\
(4) \quad & \frac{\Gamma' = \Gamma[v \leftarrow \{(\text{alloc}(T)@p, \top)\}]}{\Upsilon, \Gamma \vdash v = \text{new}^p T : \Gamma'} \\
(5) \quad & \frac{\Upsilon, \Gamma \vdash I_1 : \Gamma_1 \quad \Upsilon, \Gamma \vdash I_2 : \Gamma_2}{\Upsilon, \Gamma \vdash \text{if}(\ast) I_1 \text{ else } I_2 : \Gamma_1 \sqcup \Gamma_2} \\
(6) \quad & \frac{\Upsilon, \Gamma \vdash I_1 : \Gamma_1 \quad \Upsilon, \Gamma_1 \vdash I_2 : \Gamma_2}{\Upsilon, \Gamma \vdash I_1; I_2 : \Gamma_2}
\end{aligned}$$

Figure 8: Rules for intraprocedural analysis

4.4 Intraprocedural Analysis

Figure 8 describes the intraprocedural analysis using judgements of the form $\Upsilon, \Gamma \vdash I : \Gamma'$. which indicates that, if statement I is executed in an environment that satisfies summary environment Υ and abstract heap Γ , we obtain a new heap Γ' . Since the analysis is (partially) flow-sensitive, we distinguish between heaps Γ, Γ' before and after executing I .

Rule (1) in Figure 8 describes the analysis of assignments. Although our analysis only performs weak updates to heap objects, it does apply strong updates to variables. Hence, rule (1) updates the points-to set for (v_1, ϵ) to be $\Gamma[v_2, \epsilon]$, where the lookup operation is defined in Section 4.3. Rule (2) concerns loads of the form $v_1 = v_2.f$. Here, we first look up the points-to set θ of v_2 and then use $\Gamma[\theta, f]$ to retrieve the targets of memory locations in θ . Finally, since our analysis applies strong updates to variables, we override v_1 's existing targets and change its points-to set to $\Gamma[\theta, f]$.

Rule (3) analyzes stores $v_1.f = v_2$. First, we look up the points-to sets θ_1 and θ_2 of v_1 and v_2 . Now, the store operation will update every location o_i such that $(o_i, \phi_i) \in \theta_1$. However, since we apply only weak updates to heap objects, we preserve the existing points-to targets $\Gamma(o_i, f)$ for each o_i . Furthermore, since v_1 points to o_i under constraint ϕ_i , o_i points to elements in θ_2 only when ϕ_i holds. Hence, the new points-to set for o_i is given by $(\theta_2 \downarrow \phi_i) \sqcup \Gamma(o_i, f)$ where the \downarrow operation is given by Definition 9.

Rule (4) for memory allocations $v = \text{new}^p T$ introduces a new abstract location named $\text{alloc}(T)@p$ and assigns v_1 to this singleton. Rules (5) and (6) describe the analysis of if statements and sequencing and are fairly standard.

4.5 Method Calls and Summary Instantiation

We now describe the instantiation of summaries at call sites. Since a key part of summary instantiation is constructing the mapping from locations in the summary to those at the call

$$\begin{aligned}
& \frac{}{\mathcal{M}, \Gamma, \rho \vdash \text{inst_loc}(a_i) : \{\mathcal{M}(a_i), \top\}} \\
& \frac{\mathcal{M}, \Gamma, \rho \vdash \text{inst_loc}(\pi) : \theta}{\mathcal{M}, \Gamma, \rho \vdash \text{inst_loc}(\pi.f) : \Gamma[\theta, f]} \\
& \frac{\mathcal{M}, \Gamma, \rho \vdash \text{inst_loc}(\pi) : \theta_0 \quad \theta_i = \bigsqcup_{1 \leq j \leq n} \Gamma[\theta_{i-1}, f_j]}{\mathcal{M}, \Gamma, \rho \vdash \text{inst_loc}(\pi.(f_1 \dots f_n)^*) : \bigsqcup_{i \geq 0} \theta_i} \\
& \frac{\vec{\rho}_{\text{new}} = \text{new_ctx}(\rho, \vec{\rho})}{\mathcal{M}, \Gamma, \rho \vdash \text{inst_loc}(\text{alloc}(T)@\vec{\rho}) : \{(\text{alloc}(T)@\vec{\rho}_{\text{new}}, \top)\}} \\
& \frac{\vec{\rho}_{\text{new}} = \text{new_ctx}(\rho, \vec{\rho})}{\mathcal{M}, \Gamma, \rho \vdash \text{inst_loc}(v@\vec{\rho}) : \{(v@\vec{\rho}_{\text{new}}, \top)\}}
\end{aligned}$$

Figure 9: Rules for instantiating memory locations

site, we first describe the instantiation of memory locations, which corresponds to the mapping ζ^{-1} from Section 3.

Figure 9 describes the instantiation of locations using judgements of the form $\mathcal{M}, \Gamma, \rho \vdash \text{inst_loc}(\pi) : \theta$. Here, \mathcal{M} maps formals to actuals, and Γ and ρ are the abstract heap and program point associated with a call site respectively. The meaning of the judgement is that, under $\mathcal{M}, \Gamma, \rho$, location π used in the summary maps to (guarded) location set θ .

The first rule Figure 9 maps formal parameter a_i to the actual $\mathcal{M}(a_i)$. The second rule instantiates argument-derived locations of the form $\pi.f$. For this purpose, we first instantiate prefix π to location set θ , then retrieve the points-to targets of the f field of locations in θ . The third rule describes the instantiation of access paths of the form $\pi.(f_1 \dots f_n)^*$. As in the previous rule, we first instantiate prefix π , which yields θ_0 . Now, recall that the access path $\pi.(f_1 \dots f_n)^*$ describes the infinite set of access paths given by the regular expression $\pi.(f_1 + \dots + f_n)^*$. Hence, to instantiate $\pi.(f_1 \dots f_n)^*$, we need to compute all locations that are reachable from θ_0 using any combination of field selectors f_1, \dots, f_n . The resulting set $\bigsqcup_{i \geq 0} \theta_i$ is therefore the reflexive transitive closure of θ_0 with respect to fields f_1, \dots, f_n .

The last two rules in Figure 9 describe the instantiation of allocations and local variables. Both rules use a helper `new_ctx` method defined as follows:

$$\text{new_ctx}(\rho, \vec{\rho}) = \begin{cases} \rho, \vec{\rho} & \text{if } |\vec{\rho}| \leq k \\ \vec{\rho} & \text{otherwise} \end{cases}$$

In other words, `new_ctx` appends call site ρ to context $\vec{\rho}$ if the length of $\vec{\rho}$ is less than some pre-determined threshold k . Hence, our analysis uses a k -CFA style context-sensitive heap abstraction where the value of k is configurable.

EXAMPLE 5. Consider the following abstract heap Γ :

$$\begin{aligned}
\Gamma(x, \epsilon) &= \{(o_1, \top)\} & \Gamma(x.\epsilon, f) &= \{(o_2, \phi_1), (o_3, \phi_2)\} \\
\Gamma(o_2, f) &= \{(o_4, \top)\} & \Gamma(o_3, f) &= \{(o_5, \top)\}
\end{aligned}$$

At a call site where a_0 maps to actual x , the access path $a_0.\epsilon.f^*$ instantiates to $\{(o_2, \phi_1), (o_3, \phi_2), (o_4, \phi_1), (o_5, \phi_2)\}$.

Figure 10 describes the instantiation of constraints. To translate a constraint $\text{type}(t) = T$, we map t to its corresponding location set θ by using `inst_loc`. The function

$$\frac{\mathcal{M}, \Gamma, \rho \vdash \text{inst_loc}(\text{lift}^{-1}(t)) : \theta \quad \phi = \text{has_type}(\theta, T)}{\mathcal{M}, \Gamma, \rho \vdash \text{inst}_\phi(\text{type}(t) = T) : \phi}$$

$$\frac{\star \in \{\wedge, \vee\} \quad \mathcal{M}, \Gamma, \rho \vdash \text{inst}_\phi(\phi_1) : \phi'_1 \quad \mathcal{M}, \Gamma, \rho \vdash \text{inst}_\phi(\phi_2) : \phi'_2}{\mathcal{M}, \Gamma, \rho \vdash \text{inst}_\phi(\phi_1 \star \phi_2) : \phi'_1 \star \phi'_2}$$

Figure 10: Rules for instantiating constraints

$$\frac{\mathcal{M}, \Gamma, \rho \vdash \text{inst_loc}(\pi_1) : \theta_1 \dots \text{inst_loc}(\pi_n) : \theta_n \quad \mathcal{M}, \Gamma, \rho \vdash \text{inst}_\phi(\phi_1) : \phi'_1 \dots \text{inst}_\phi(\phi_n) : \phi'_n}{\mathcal{M}, \Gamma, \rho \vdash \text{inst_pts}(\{(\pi_1, \phi_1), \dots, (\pi_n, \phi_n)\}) : \sqcup_i(\theta_i \downarrow \phi_i)}$$

$$\frac{\mathcal{M}, \Gamma, \rho \vdash \text{inst_loc}(\pi) : \theta' \quad \mathcal{M}, \Gamma, \rho \vdash \text{inst_pts}(\theta) : \theta'' \quad \Delta = [(\pi_i, f) \leftarrow (\theta'' \downarrow \phi_i) \mid (\pi_i, \phi_i) \in \theta']}{\mathcal{M}, \Gamma, \rho \vdash \text{inst_partial_heap}(\pi, f, \theta) : \Delta}$$

$$\frac{\Delta = \{(\pi_1, f_{11}) \mapsto \theta_{11}, \dots, (\pi_n, f_{nk}) \mapsto \theta_{nk}\} \quad \mathcal{M}, \Gamma, \rho \vdash \text{inst_partial_heap}(\pi_1, f_{11}, \theta_{11}) : \Delta_{11} \quad \dots \quad \mathcal{M}, \Gamma, \rho \vdash \text{inst_partial_heap}(\pi_n, f_{nk}, \theta_{nk}) : \Delta_{nk}}{\mathcal{M}, \Gamma, \rho \vdash \text{inst_heap}(\Delta) : \sqcup_{ij} \Delta_{ij}}$$

Figure 11: Rules for instantiating summaries

$\text{has_type}(\theta, T)$ then yields the condition under which some element in θ has dynamic type T (recall Definition 5).

Figure 11 describes the instantiation of heap Δ . Given location π_i and field f_j from the callee heap, inst_partial_heap instantiates all points-to edges from π_i labeled with f_j and yields instantiated partial heap Δ_{ij} . The instantiation of Δ is obtained by taking the join over all Δ_{ij} ’s.

Finally, Figure 12 describes the analysis of method calls. First, consider a static call to m with corresponding summary Δ (rule (1)). To analyze it, we construct the formal-to-actual mapping \mathcal{M} and perform a least fixed-point computation that instantiates the summarized heap Δ until we obtain an overapproximation of the set Δ' of m ’s side effects. The abstract heap after the method call is obtained by taking the union of the existing heap Γ and the new “edges” Δ' .

The second rule of Figure 12 describes the analysis of virtual calls. Here, we first overapproximate the call’s targets and then use the previous rule for analyzing static method calls to obtain heap Γ_i assuming the called method is $T_i :: m$. Now, since the target of the virtual call is $T_i :: m$ under the assumption that v_0 has dynamic type T_i , we generate the constraint $\phi_i = \text{has_type}(\Gamma(v_0), T_i)$. The final abstract heap after the call is obtained as $\sqcup_i(\Gamma_i \downarrow \phi_i)$ where the notation $\Gamma_i \downarrow \phi_i$ conjoins ϕ_i with all the constraints in Γ_i .

4.6 Summary Generation

We now consider the last piece of the analysis, which is summary generation. As shown in Figure 13, the summary Δ for a method m is obtained by analyzing m ’s body I starting with the *empty abstract heap*. In particular, the initial heap does not contain any nodes or edges because our analysis constructs the normal form heap on demand depending on which fields are accessed in I . Furthermore, since we do

$$(1) \quad \frac{\Upsilon(T, m) = \Delta \quad \mathcal{M} = [a_1 \mapsto v_1, \dots, a_n \mapsto v_n] \quad \mathcal{M}, \Gamma \sqcup \Delta', \rho \vdash \text{inst_heap}(\Delta) : \Delta'}{\Upsilon, \Gamma \vdash m^\rho @ T(v_1, \dots, v_n) : \Gamma \sqcup \Delta'}$$

$$\frac{\text{static_type}(v_0) = T \quad T_1 <: T, \dots, T_n <: T \quad \phi_i = \text{has_type}(\Gamma(v_0), T_i) \quad \Upsilon, \Gamma \vdash m^\rho @ T_1(v_0, \dots, v_k) : \Gamma_1 \quad \dots \quad \Upsilon, \Gamma \vdash m^\rho @ T_n(v_0, \dots, v_k) : \Gamma_n}{\Upsilon, \Gamma \vdash v_0.m^\rho(v_1, \dots, v_k) : \sqcup_i(\Gamma_i \downarrow \phi_i)}$$

Figure 12: Analysis of method calls

$$\frac{\Upsilon, \emptyset \vdash I : \Delta \quad \Upsilon(T, m) = \Delta}{\Upsilon \vdash T :: m(T_o v_0, \dots, T_n v_n) \{V^*; I\}}$$

Figure 13: Summary generation

not explicitly add default edges, the resulting heap Δ only contains side effects introduced by method m .

Note that the summary for the program’s entry method (i.e., main) can be used to determine the points-to sets for all program variables. If the program does not have a unique main method, then the final points-to sets are obtained by taking the join of the summaries of all entry methods.

5. IMPLEMENTATION AND EXTENSIONS

We implemented the proposed algorithm in a tool called SCUBA which consists of $\sim 19,000$ lines of Java code. Our implementation is built on top of Chord [15] and performs analysis on the Quad intermediate representation of Joeq [25]. SCUBA obtains an initial callgraph by running the context-insensitive pointer analysis implemented in Chord. Our implementation also uses the Z3 SMT solver [5] for checking satisfiability of constraints.

Our implementation performs several optimizations over the core algorithm described here. One optimization is memoizing instantiation results. For example, consider an access path $a_0.e.f.g.h$ used in m ’s summary. Since this access path may be instantiated many times when analyzing a call site of m , our analysis maintains a cache per callsite that records instantiation results. A second optimization concerns constraint generation for virtual method calls. Consider a call $v.m(\dots)$ where the static type of v is T_0 . Further, suppose T_0 has a large number of subclasses T_1, \dots, T_n all of which inherit T_0 ’s m method except T_n . Assuming v points to heap object o , we need to introduce constraints of the form $\bigvee_{0 \leq i < n} \text{type}(o) = T_i$. Since such constraints can be very large, our implementation allows subtyping constraints and translates them to linear inequalities by assigning integer identifiers to types in reverse topological order.

Another optimization concerns the tracking of global variables. In particular, in addition to maintaining a per-method *local method summary*, we also maintain a *global summary*, which includes information about side-effects to globals. Since we use this global summary to retrieve points-to targets of static fields, we do not need to instantiate a large number of global variables at every call site. Finally, we resort to a

Benchmark	# methods	# statements	CIPA	2-CFA	2-obj	SCUBA-2	SCUBA-3	SCUBA-4
antlr	5411	112831	29	1380	355	30	37	50
hedc	4967	103066	25	1337	446	25	28	31
avroa	5230	104948	24	1328	336	53	55	59
polyglot	4634	92615	21	608	284	14	17	16
toba-s	4702	101501	24	930	299	15	18	17
weblech	5816	115937	27	1657	506	41	35	39
xalan	6405	131332	27	1100	3600	173	180	211
hsqldb	6767	137947	33	2474	1348	63	65	77
luindex	6157	127451	28	2525	532	93	147	315
sunflow	9653	227507	66	T/O	T/O	411	405	521
Average	5974	125514	30	1694	1131	92	99	134

Table 1: Analysis running time in seconds. Each analysis is run with a time limit of 3600 seconds; runs exceeding this limit are labeled T/O. For analyses that time out, we use the lower bound of 3600 s. for computing averages.

number of heuristics while analyzing SCCs in the callgraph. For example, since the size of constraints can blow up during analysis of SCCs, we over-approximate constraints when their sizes grow beyond a pre-determined threshold.

6. EVALUATION

We evaluated our algorithm by running SCUBA on ten large Java applications from the DaCapo and Ashes benchmark suites [2, 1]. These applications range between 92615 and 227507 lines of statements in the Quad IR and contain between 4634 and 9653 reachable methods.

To evaluate our algorithm, we compared SCUBA against the context-insensitive as well as k -CFA and k -object-sensitive algorithms implemented in Chord [15]. All analyses implemented in Chord are Anderson-style flow-insensitive pointer analyses that allow customizing the value of k . Chord also allows customizing the context-sensitivity associated with heap objects using a value h . For example, a 2-obj-1- h analysis uses the abstract allocation site of the receiver as a context up to depth 2, and it also differentiates heap allocations with different contexts up to depth 2.

Before we describe our experimental results, we first explain how our algorithm conceptually relates to k -CFA and k -object-sensitive analyses. Similar to k -CFA, SCUBA uses *call sites* rather than *receiver objects* as contexts. However, unlike k -CFA, we do not impose a fixed value of k since our algorithm instantiates method summaries differently for each call site. On the other hand, we can customize the context-sensitivity associated with heap objects by varying the parameter k used in the `new_ctx` function from Section 4.5. Hence, for a given value of k in the `new_ctx` function, SCUBA is roughly comparable to a ∞ -CFA- k - h analysis in the terminology of top-down pointer analysis. In what follows, we use the notation SCUBA- k to refer to different configurations of SCUBA for different values of parameter k used in the `new_ctx` function from Section 4.5.

Table 1 compares the running times of SCUBA- k (for $2 \leq k \leq 4$) against the context-insensitive, 2-CFA, and 2-object-sensitive analyses using h value of 1. While we also tried comparing SCUBA against k -CFA and k -object-sensitive analyses for $k = 3$ and $k = 4$, these analyses did not complete within an hour for *most of the benchmarks*; hence, we do not include these results in Table 1. We also note that the running times shown in Figure 1 include the analysis time for libraries (e.g., JDK, Swing, Sun Security) as well as the application code. In other words, we did not use manually

provided stub methods for analyzing libraries. As the results in Table 1 show, SCUBA- k is significantly faster compared to k -CFA and k -obj analyses. For example, for $k = 2$, SCUBA is approximately $12\times$ faster than the object-sensitive analysis and about $18\times$ faster than the k -CFA analysis on average.

To compare the precision of SCUBA against k -CFA and k -obj analyses, we used two typical pointer analysis clients, namely *may-alias* and *downcast* analyses. Table 2 compares the precision of different analysis configurations in the context of the may-alias client. The column labeled “Alias pairs” shows the number of variable pairs that are queried by the may-alias client. To generate these pairs, we first ran a context-insensitive pointer analysis to identify potential may-alias in the application code. From these variables, we further filtered those pairs that are “obviously” aliases (e.g., due to a direct assignment).

Columns 3-11 in Table 2 show the number of variables proven not to be aliases according to each analysis configuration. Hence, a higher number indicates better precision. Observe that every configuration of SCUBA- k yields better precision *on average* compared to 2-CFA and 2-obj analyses ⁴.

Table 3 shows the precision of each analysis in the context of the downcast client. Here, the second column labeled “# downcasts” shows the total number of downcasts in the application, and the subsequent columns show the number of downcasts that can be proven safe. The number of downcasts shown in Table 3 only include the downcasts performed in the application code rather than in external libraries ⁵. According to the results shown in Table 3, SCUBA- k has better precision on average compared to both 2-CFA and 2-obj.

7. RELATED WORK

Top-down pointer analysis. Most existing context-sensitive pointer analysis algorithms are top-down [11, 13, 26, 24, 7]. Generally speaking, top-down context-sensitivity comes in two flavors: *call-site sensitivity* [22] (k -CFA) and *k-object-sensitivity* [13]. Specifically, CFA-based algorithms use method call sites as the context, while object-sensitive approaches use the receiver’s abstract allocation site. The recent work

⁴We manually inspected a randomly selected subset of the may-alias queries that could only be discharged by SCUBA and confirmed that these are not false negatives.

⁵Since most benchmarks use the same libraries, this strategy avoids double counting the results. Furthermore, bug finding and verification clients are typically interested in finding defects in the application (rather than library) code.

Benchmark	Alias pairs	CIPA	2-CFA	2-obj	SCUBA-2	SCUBA-3	SCUBA-4
antlr	6839	0	1082	2785	3219	3231	3231
hedc	1728	0	725	962	1025	1055	1055
avroara	1182	0	406	687	738	741	745
polyglot	165	0	59	103	128	128	128
toba-s	5118	0	3354	3350	3589	3589	3595
weblech	1417	0	662	654	656	681	763
xalan	124	0	24	24	24	24	24
hsqldb	5254	0	2746	2724	3318	3318	3426
luindex	4649	0	1326	1420	1353	1400	1400
sunflow	4303	0	N/A	N/A	339	339	339
Average	2942	0	1154	1412	1561	1574	1596

Table 2: May-alias results. Bigger numbers indicate higher precision; row labeled “Average” excludes sunflow.

Benchmark	# downcasts	CIPA	2-CFA	2-obj	SCUBA-2	SCUBA-3	SCUBA-4
antlr	76	18	21	48	45	52	52
hedc	28	5	6	23	23	23	23
avroara	21	0	0	8	15	18	18
polyglot	13	2	3	9	13	13	13
toba-s	59	23	23	34	37	37	37
weblech	48	16	23	33	38	38	38
xalan	14	7	10	13	13	13	13
hsqldb	45	22	24	32	28	28	35
luindex	213	104	106	180	177	177	177
sunflow	81	19	N/A	N/A	25	52	52
Average	57	22	24	42	43	44	45

Table 3: Downcast results. Bigger numbers indicate higher precision; row labeled “Average” excludes sunflow.

described in [8] has proposed selectively combining k -CFA and k -object sensitivity to achieve better precision.

Several papers have used BDD-based methods for top-down context-sensitive pointer analysis [26, 31, 10]. The use of BDDs exposes commonalities between different contexts and allows the technique to scale better. The Chord framework [15] used in our experimental evaluation also uses BDDs to exploit equivalences between calling contexts.

Bottom-up pointer analysis. While not as widely-studied as top-down algorithms, several papers propose bottom-up pointer analysis. However, many of these approaches are unification-based [9, 17, 19, 12]. The algorithm proposed in [28] is also bottom-up but uses a combined equality- and subset-based approach. By contrast, our algorithm is subset-based and therefore more precise.

The algorithm described in [4] presents a subset-based, partially bottom-up pointer analysis for C. Unlike our approach, this algorithm incorporates both top-down and bottom-up phases and applies these analyses until a fixed-point is reached. Here, the top-down phase is used for precise handling of function pointers, while the bottom-up phase is used for propagation of points-to relations. Furthermore, the algorithm described in [4] tracks alias pairs as opposed to an explicit heap model and is meant for C rather than Java. Another subset-based pointer analysis for C that combines top-down and bottom-up phases is described in [16]. This algorithm is based on the observation that context-insensitivity does not result in a loss of precision if function side effects are accounted for. In contrast to our algorithm, [16] handles SCCs in a completely context-insensitive way and employs a top-down phase that removes callee side effects.

The algorithms described in [3, 6] perform bottom-up pointer

analysis for C++ programs. Unlike the method presented in this paper, they perform case splits on possible aliasing patterns. Since the approach of [6] performs strong updates, it is more precise but less scalable compared to our technique.

Another related work is the algorithm described in [27], which describes a compositional pointer and escape analysis for Java. While this analysis is flow-sensitive and applies strong updates, it assumes that parameters do not alias and generates a summary that is valid under this assumption. However, if this assumption is violated at a call site, the analysis corrects the summary through a complex mechanism that involves merging of memory locations. Another difference is that [27] does not precisely handle virtual method calls and uses a CHA-based callgraph.

Summary-based interprocedural analysis. Many papers describe general frameworks for interprocedural analysis [18, 20, 21]. More recently, the work described in [29] and [30] compute polymorphic summaries for dataflow problems but both rely on global points-to sets.

8. CONCLUSION

We described a new bottom-up, summary-based pointer analysis for Java. The experimental evaluation demonstrates that our algorithm runs significantly faster than top-down pointer analyses with comparable precision. We believe that SCUBA is able to scale better because the cost of instantiating a method summary is smaller compared to the cost of re-analyzing the function.

In future work, we plan to use the proposed analysis in the context of more bug finding and verification clients (e.g., tpestate analysis). We also to plan to adapt the techniques described in this paper for demand-driven bottom-up pointer analysis.

9. REFERENCES

- [1] Ashes benchmark suite.
<http://www.sable.mcgill.ca/software/#ashessuitecollection>.
- [2] Dacapo benchmarks. <http://www.dacapobench.org/>.
- [3] Ramkrishna Chatterjee, Barbara G. Ryder, and William A. Landi. Relevant context inference, 1999.
- [4] Ben-Chung Cheng and Wen-Mei W. Hwu. Modular interprocedural pointer analysis using access paths: Design, implementation, and evaluation. *PLDI '00*, pages 57–69, 2000.
- [5] Leonardo De Moura and Nikolaj Bjørner. Z3: An efficient smt solver. In *TACAS*, pages 337–340. Springer, 2008.
- [6] Isil Dillig, Thomas Dillig, Alex Aiken, and Mooly Sagiv. Precise and compact modular procedure summaries for heap manipulating programs. In *PLDI*, pages 567–577, 2011.
- [7] Manuel Fähndrich, Jakob Rehof, and Manuvir Das. Scalable context-sensitive flow analysis using instantiation constraints. *PLDI'00*, 2000.
- [8] George Kastrinis and Yannis Smaragdakis. Hybrid context-sensitivity for points-to analysis. *PLDI*, pages 423–434, 2013.
- [9] Chris Lattner, Andrew Lenharth, and Vikram Adve. Making context-sensitive points-to analysis with heap cloning practical for the real world. In *PLDI*, 2007.
- [10] Ondrej Lhoták. *Program analysis using binary decision diagrams*. PhD thesis, McGill University, 2006.
- [11] Ondřej Lhoták and Laurie Hendren. Context-sensitive points-to analysis: is it worth it? In *CC*, pages 47–64, 2006.
- [12] Donglin Liang and Mary Jean Harrold. Efficient points-to analysis for whole-program analysis. In *FSE*, pages 199–215, 1999.
- [13] Ana Milanova, Atanas Rountev, and Barbara G. Ryder. Parameterized object sensitivity for points-to and side-effect analyses for java. In *ISSTA*, pages 1–11, New York, NY, USA, 2002.
- [14] Ana Milanova, Atanas Rountev, and Barbara G. Ryder. Parameterized object sensitivity for points-to analysis for java. *ACM Trans. Softw. Eng. Methodol.*, 2005.
- [15] Mayur Naik. Chord framework.
<http://pag.gatech.edu/chord>.
- [16] Erik M Nystrom, Hong-Seok Kim, and W Hwu Wen-mei. Bottom-up and top-down context-sensitive summary-based pointer analysis. In *Static Analysis*, pages 165–180. Springer, 2004.
- [17] Robert O’Callahan. *Generalized aliasing as a basis for program analysis tools*. PhD thesis, Carnegie Mellon University, 2001.
- [18] Thomas Reps, Susan Horwitz, and Mooly Sagiv. Precise interprocedural dataflow analysis via graph reachability. *POPL '95*, pages 49–61.
- [19] Erik Ruf. Effective synchronization removal for java. In *ACM SIGPLAN Notices*, volume 35, pages 208–218. ACM, 2000.
- [20] Shmuel Sagiv, Thomas W. Reps, and Susan Horwitz. Precise interprocedural dataflow analysis with applications to constant propagation. *Theor. Comput. Sci.*, 167(1&2):131–170, 1996.
- [21] Micha Sharir and Amir Pnueli. *Two approaches to interprocedural data flow analysis*, chapter 7, pages 189–234. Prentice-Hall, Englewood Cliffs, NJ, 1981.
- [22] Olin Shivers. Control-flow analysis of higher-order languages. Technical report, 1991.
- [23] Yannis Smaragdakis, Martin Bravenboer, and Ondrej Lhoták. Pick your contexts well: Understanding object-sensitivity. *POPL*, 2011.
- [24] Manu Sridharan and Rastislav Bodík. Refinement-based context-sensitive points-to analysis for java. *PLDI*, pages 387–400, 2006.
- [25] John Whaley. Joeq: A virtual machine and compiler infrastructure. In *Proceedings of the 2003 workshop on Interpreters, virtual machines and emulators*, pages 58–66. ACM, 2003.
- [26] John Whaley and Monica S Lam. Cloning-based context-sensitive pointer alias analysis using binary decision diagrams. In *PLDI*, pages 131–144, 2004.
- [27] John Whaley and Martin Rinard. Compositional pointer and escape analysis for java programs. *OOPSLA*, pages 187–206, 1999.
- [28] Guoqing Xu and Atanas Rountev. Merging equivalent contexts for scalable heap-cloning-based context-sensitive points-to analysis. In *Proceedings of the 2008 international symposium on Software testing and analysis*, pages 225–236. ACM, 2008.
- [29] Greta Yorsh, Eran Yahav, and Satish Chandra. Generating precise and concise procedure summaries. In *POPL*, volume 43, pages 221–234, 2008.
- [30] Xin Zhang, Ravi Mangal, Mayur Naik, and Hongseok Yang. Hybrid top-down and bottom-up interprocedural analysis. In *PLDI*, page 28, 2014.
- [31] Jianwen Zhu and Silvan Calman. Symbolic pointer analysis revisited. In *PLDI*, pages 145–157, 2004.