

Precise Identification of Composition Relationships for UML Class Diagrams

Ana Milanova

Department of Computer Science
Rensselaer Polytechnic Institute
milanova@cs.rpi.edu

Abstract

Knowing which associations are compositions is important in a tool for the reverse engineering of UML class diagrams. Firstly, recovery of composition relationship bridges the gap between design and code. Secondly, since composition relationships explicitly state a requirement that certain representation cannot be exposed, it is important to determine if this requirement is met by component code. Verifying that compositions are implemented properly may prevent serious program flaws due to representation exposure.

We propose an implementation-level composition model based on ownership and a novel approach for identifying compositions in Java software. Our approach uses static ownership inference based on points-to analysis and is designed to work on incomplete programs. We present empirical results on several components. In our experiments, on average 40% of the examined fields account for relationships that are identified as compositions. We also present a precision evaluation which shows that our analysis achieves almost perfect precision—that is, it almost never misses composition relationships. The results indicate that precise identification of interclass relationships can be done with a simple and inexpensive analysis, and thus can be easily incorporated in reverse engineering tools that support iterative model-driven development.

1. Introduction

In modern software development design recovery through reverse engineering is performed often; in a typical iterative development process reverse engineering is performed at the beginning of every iteration to recover the design from the previous iteration [11].

UML class diagrams describe the architecture of the program in terms of classes and interclass relationships; they are scalable, informative and widely-used design models. While the UML concepts of class and inheritance have corresponding first-class concepts in object-oriented programming languages, the UML concepts of *association*, *aggre-*

gation and *composition* do not have corresponding language concepts. Thus, while the reverse engineering of classes and inheritance hierarchies is straightforward, the reverse engineering of associations presents various challenges.

UML associations model relatively permanent interclass relationships; conventionally, they are implemented using instance fields of reference type [11] (e.g., an association from class *A* to class *B* is implemented using a reference field of type *B* in class *A*). Thus, reverse engineering tools infer associations by examining instance fields of reference type; however, the inference is often non-trivial. One challenge is the recovery of one-to-many associations implemented using pseudo-generic containers (e.g., `Vector`). Another challenge is the recovery of compositions. Modern reverse engineering tools such as Rational ROSE and ArgoUML do not address these challenges and produce inconsistent and even incorrect class diagrams (see Guéhéneuc and Albin-Amiot [9] for detailed examples). Clearly, this leads to a gap between design class diagrams and reverse engineered class diagrams which hinders understanding, round-trip engineering and identification of design patterns.

Towards the goal of bridging this gap, this paper proposes a methodology for inference of binary associations for UML class diagrams. Its major emphasis and contribution is the inference of composition relationships, which we believe is challenging and inadequately addressed in previous work. While the UML concept of aggregation is “largely meaningless” [6], the UML concept of composition has a well-defined semantics that emphasizes the notion of *ownership*: a “composition is a strong form of [whole-part] association with strong ownership of parts by the composite and coincident lifetime of parts with the composite. A part may belong to only one component at a time” [18, Chapter 14]. Therefore, a composition relationship at design level states the requirement for ownership and no *representation exposure* at implementation level (i.e., the owned component object cannot be exposed outside of its composite owner object); if composition is implemented properly ownership should be preserved.

It is important to investigate techniques for recovery of

composition relationships. Firstly, it helps bridge the gap between the design class diagram and the reverse engineered diagram. Secondly, since composition relationships explicitly state a requirement that certain representation cannot be exposed, it is important to determine if this requirement is met by component code. Verifying that compositions are implemented properly may prevent serious program flaws due to representation exposure such as the well-known `signers` bug in Java 1.1.¹

Thus, the goals of this work are (i) to define an implementation-level ownership model that captures the notion of composition in design and (ii) to design an analysis algorithm that infers ownership and composition in accordance with this model. Our definition of implementation-level composition is based on the *owners-as-dominators* ownership model [4, 13]; in this model the owner object (the composite) should dominate an owned object (a component)—that is, all access paths to the owned object should pass through its owner. The owners-as-dominators model defines an ownership boundary for each owner; intuitively, an owned object may be accessed by its owner as well as other objects within the boundary of the owner. For example, an owned object stored in an instance field may be passed to an owned container. As pointed out by Clarke et. al [4, 13] and also observed by us during the empirical investigation, the owners-as-dominators model captures well the notion of composition in modeling.

We propose a novel static analysis for ownership inference. If the ownership inference determines that all objects stored in a field are owned by their enclosing object, the analysis identifies a composition through that field. Our approach works on incomplete programs. This is an important feature because in the context of reverse engineering tools it is essential to be able to perform separate analysis of software components. For example, it is typical to have to analyze a component without having access to the clients of that component. Our ownership inference analysis is based on *points-to analysis*, which determines the set of objects a reference variable or a reference object field may point to. We use the points-to analysis solution to approximate the possible accesses between run-time objects.

We present empirical results on several components. In our experiments, on average 40% of the examined fields account for relationships that are identified as compositions. We also present a precision evaluation which shows that our analysis achieves almost perfect precision—that is, it almost never misses composition relationships identified in our model. The results indicate that precise identification of interclass relationships can be done with a simple and inex-

pensive analysis, and thus can be easily incorporated in reverse engineering tools that support iterative model-driven development.

This work has the following contributions:

- We propose an implementation-level ownership and composition model that captures well the notion of composition in modeling.
- We propose a static analysis for identifying composition relationships in accordance with our model; the analysis works on incomplete programs.
- We present an empirical study that evaluates our analysis on several Java components.

2. Problem Statement

Reverse engineering tools typically infer associations by examining instance fields of reference type in the code. In our model, an association relationship through a field f is refined as composition if it can be proven that all objects referred by f are owned by their enclosing object. Thus, given a suitable definition of implementation-level ownership and composition, our goal is to design a static analysis that answers the question: given a set of Java classes (i.e, a component to be analyzed) for what instance fields we observe implementation-level composition throughout all possible executions of arbitrary client code built on top of these classes? The analysis output is a set of fields for which the relationship is guaranteed to be a composition for arbitrary client code.

Clearly, there are various uses of this information when integrated in a tool for the reverse engineering of UML class diagrams. Firstly, recovery of compositions bridges the gap between design and implementation and eases the understanding of underlying design. Secondly, verifying that compositions are implemented properly and refactoring early if needed may help prevent serious program flaws such as the `signers` bug in Java 1.1. Ownership of representation is a desirable property and information about it can be included automatically in the component documentation.

The input to the analysis contains a set Cls of interacting Java classes. We will use "classes" to denote both Java classes and interfaces because for the purposes of this work the difference is not relevant. A subset of Cls is designated as the set of *accessible classes*; these are classes that may be accessed by unknown client code from outside of Cls . Such client code can only access fields and methods from Cls that are declared in some accessible class; these accessible fields and methods are referred to as *boundary fields* and *boundary methods*.

Section 2.1 describes the ownership model, and Section 2.2 describes the notion of implementation-level com-

¹ In Java 1.1 the security system function `Class.getSigners` returned a pointer to an internal array allowing clients to modify the array and compromising the security of the system.

```

public class Vector {
  protected Object[] data;
  public Vector(int size) {
1 data = new Object[size]; }
  public void addElement(Object e,int at) {
2 data[at] = e; }
  public Object elementAt(int at) {
3 return data[at]; }
  public Enumeration elements() {
4 return new VIterator(this); } }

final class VIterator implements Enumeration {
  Vector vector;
  int count;
  VIterator(Vector v) {
5 this.vector = v;
6 this.count = 0; }
  Object nextElement() {
7 Object[] data = vector.data;
8 int i = this.count;
9 this.count++;
10 return data[i]; } }

main() {
11 Vector v = new Vector(100);
12 X x = new X();
13 v.addElement(x,0);
14 Enumeration e = v.elements();
15 x = (X) e.nextElement();
16 x.m(); }

```

Figure 1. Simplified vector and its iterator.

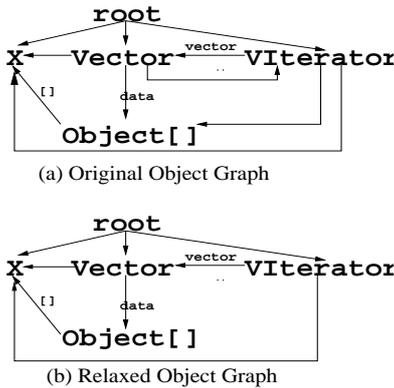


Figure 2. Object graphs for Figure 1.

position based on it. Section 2.3 discusses certain constraints to the model that allow more precise identification of ownership and composition.

2.1. Ownership Model

The ownership model is based on the notion of *owners-as-dominators* [4, 3, 13]. It is essentially the model proposed by Potter et al. [13] with several modifications that allow more precise handling of popular object-oriented patterns and idioms such as iterators, composites and factories [7]. In this model, each execution is represented by an

object graph which describes access relationships between run-time objects:

- Let f be a reference instance field in a run-time object o . There is an edge $o \xrightarrow{f} o'$ in the object graph iff field f in o refers to o' at some point of program execution.²
- There is an edge $o \Downarrow o'$ iff some element of array o refers to o' at some point of program execution.
- There is an edge $o \rightarrow o'$ iff an instance method or constructor invoked on receiver o has local variable r that refers to o' , or a static method called from an instance method or constructor invoked on o , has a local variable r that refers to o' . There is an edge of this kind only if there is no edge of the first kind from o to o' .

A run-time object o' is accessed in the *context* of o iff there is an edge from o to o' in the object graph. The start of program execution is expressed with a special node *root*. Context *root* represents the context for *main* and for objects referenced by static fields. For example, executing *main* in Figure 1 results in the object graph in Figure 2(a). Node *Vector* corresponds to the object created at the *new* site at line 11, node *Object[]* corresponds to the array created at the site at line 1, node *VIterator* corresponds to the iterator created at the site at line 4, and node *X* corresponds to the object created at the site at line 12.

The owners-as-dominators model states that the owner of an object o is the immediate dominator of o in the object graph [13].³ Thus, according to this model *Object[]* is not owned by its enclosing *Vector* object for this execution due to the access relationship (although only temporary) between *VIterator* and *Object[]*. To make the model less restrictive, we introduce the *relaxed object graph* which omits edges due to certain temporary access relationships. We consider two kinds of temporary access relationships. The first kind arises when an object is created in one context and immediately passed to another context; the relationship between the creating object and the new object is only temporary but if shown on the graph it is likely to restrict ownership. This notion captures the situations when an object is created and immediately returned (e.g., as in `return new VIterator(this);` in method `elements` in Figure 1) and when an object is created and immediately passed to another context (e.g., as in `new BufferedReader(new FileReader(fileName))`). This situation occurs in popular object-oriented design patterns

2 We require that all newly created objects appear in the object graph explicitly. Analogously to [4], this is done by requiring that at the point of creation a new object is stored in a new local variable; clearly, this transformation does not change program semantics.

3 Node m dominates node n if every path from the root of the graph that reaches node n has to pass through node m . The root dominates all nodes. Node m immediately dominates node n if m dominates n and there is no node p such that m dominates p and p dominates n .

such as factories, decorators and composites; in these cases the temporary relationship between the creating object and the newly created one is a matter of safety and flexibility of the implementation rather than an intention of the design. The second kind of temporary access relationships arises from field read statements $r = l.f$, where r is not assigned, passed as an implicit or explicit argument, or returned. This notion captures the situation that arises in iterators (consider statement `data = vector.data` in `nextElement` in Figure 1)—iterator objects have temporary references to the representation of their collections, which allows efficient access of collection elements; however, the collection object is always in scope. Therefore, if all accesses of o' in the context of o are due to such temporary access relationships, edge $o \rightarrow o'$ is not shown in the relaxed object graph.

The relaxed object graph for the execution of the code in Figure 1 is shown in Figure 2(b). Note that edge `Vector`→`VIterator` is omitted because it is due to a temporary access relationship of the first kind; edge `VIterator`→`Object[]` is omitted as well because it is due to a temporary access relationship of the second kind. The owner of o is the immediate dominator of o in the relaxed object graph. Thus, `root` owns `X`, `Vector` and `VIterator` and `Vector` owns `Object[]`.

2.2. Implementation-level Composition

Let A be a class in Cls , and f be a field of type B declared in A where B is a reference type (class, interface or array type [8]). The ownership property holds for f if throughout all possible executions of arbitrary clients of Cls , every instance of A owns the instances of B that its f field refers to. Consider the case when f is a collection field—that is, all objects stored in the field are arrays or instances of one of the standard `java.util` collection classes (e.g., `java.util.Vector`). If every instance of A owns all corresponding instances stored in the collection, there is a *one-to-many composition* relationship between A and C , where C is the lowest common supertype of the instances stored in the collection; otherwise, there is a one-to-many regular association (it may also be referred as aggregation). For collection fields for which the ownership property holds, there is an attribute of the association *{owned collection}* which indicates that the collection is owned by its enclosing object. Consider the case when f is not a collection field. If the ownership property holds for f , the association between A and B is a *one-to-one composition*; otherwise it is a regular one-to-one association.

Note that in Java, due to interfaces, it is not always possible to find unique non-trivial (i.e., not `Object`) least common supertype of two types. Thus, we use a variant of the Java type system in which any two types have a unique

lowest-common supertype, used elsewhere as well [20, 5]. In this system types include a powerset of the Java types in the classes being analyzed. The mapping from a Java type, T is the smallest set that includes T and all its supertypes in Java. Since the elements in this system are sets, the lowest common supertype is the intersection of two sets.

Consider the following example taken from [20]. Suppose that the least common supertype of two classes C_1 and C_2 , both of which implement interfaces I_1 and I_2 is needed. Since there are two common immediate supertypes of C_1 and C_2 , there is no unique supertype. The mapping of C_1 in the new type system is $\{C_1, I_1, I_2, \text{Object}\}$; similarly, the mapping of C_2 is $\{C_2, I_1, I_2, \text{Object}\}$. The least common supertype in these types is simply the intersection of the two sets: $\{I_1, I_2, \text{Object}\}$. When mapping the new types back to Java types it is not guaranteed that a Java type will exist. Thus, an analysis can either choose an alternative real type such as `Object`, show one-to-many associations with all types in the set, or ask the user for help. Our analysis, described in Section 4, would display `Object` in this case; assuming that the analysis is relatively precise, having `Object` as the least common supertype would indicate a potential flaw and trigger a careful review of code and design.

Example. Consider the package in Figure 3. This example is based on classes from the standard Java library package `java.util.zip`, with some modifications made to simplify the presentation and better illustrate the problem and our approach. Cls contains the classes from Figure 3 plus class `ZipEntry`. The accessible classes are `ZipInputStream`, `ZipOutputStream` and `ZipEntry` and the boundary methods are all public methods declared in those classes (i.e., the component can be accessed from client code through the public methods declared in these classes).

Clearly, the `CRC32` objects are always owned by their enclosing streams. Therefore, there is one-to-one composition relationships through fields `crc` in both `ZipInputStream` and `ZipOutputStream`. There is a regular one-to-one association through field `entry` in `ZipInputStream`; it is easy to construct client code on top of these classes such that the `ZipEntry` instances created in `ZipInputStream` objects are leaked to client code from `getNextEntry`. Similarly, there is a regular one-to-one association through `entry` in `ZipOutputStream` because the `ZipEntry` objects are passed from client code to `putNextEntry`. The associations through fields `names` and `entries` are both one-to-many regular associations between `ZipOutputStream` and `ZipEntry`; both have attribute *{owned collection}*. Clearly, the `ZipOutputStream` instance trivially owns the `Hashtable` instance. It owns the `Vector` instance as well, although the `Vector` instance is referred

```

package zip;
public class InflaterInputStream {
    protected Inflater inf;
    protected byte[] buf;
    public InflaterInputStream(Inflater inf,
        int size) {
        this.inf=inf;
        buf=new byte[size]; }
    public InflaterInputStream(Inflater inf) {
        this(inf, 512); }
    // methods read and fill contain instance calls on inf }

public class ZipInputStream extends
InflaterInputStream {
    private ZipEntry entry;
    private CRC32 crc=new CRC32();
    public ZipInputStream() {
        super(new Inflater(true), 512); }
    public ZipEntry getNextEntry() {
        crc.reset();
        inf.reset();
        if ((entry=readLOC())!=null) return null;
        return entry; }
    private ZipEntry readLOC() {
        ZipEntry e=new ZipEntry();
        // code reads and writes fields of e
        return e; } }

public class ZipOutputStream extends
DeflaterOutputStream {
    private ZipEntry entry;
    private Vector entries=new Vector();
    private Hashtable names=new Hashtable();
    private CRC32 crc=new CRC32();
    public ZipOutputStream() {
        super(new Deflater(...)); }
    public void putNextEntry(ZipEntry e) {
        // code reads and writes fields of e
        if (names.put(e.name, e)!=null) { ... }
        entries.addElement(e);
        entry=e; }
    public void closeEntry() {
        ZipEntry e=entry;
        // code reads and writes fields of e
        crc.reset();
        entry=null; }
    public void finish() {
        Enumeration enum=entries.elements();
        while (enum.hasMoreElements()) { ... } } }

```

Figure 3. Sample package zip.

to in the context of its iterator (recall the example in Figure 1); however, the iterator is a local object owned by the enclosing `ZipOutputStream` object which ensures that the `Vector` instance is dominated by the enclosing `ZipOutputStream` and may be accessed only within its ownership boundary.

2.3. Discussion

In order to allow more precise identification of implementation-level composition, we employ the following constraint, standard for other problem definitions that require analysis of incomplete programs [17, 15]. We only consider executions in which the invocation of a bound-

ary method does not leave *Cls*—that is, all of its transitive callees are also in *Cls*. In particular, if we consider the possibility of unknown subclasses, all instance calls from *Cls* could potentially be “redirected” to unknown external code that may affect the composition inference. For example, a field may be identified as composition in the current set of classes but an unknown subclass may override some method and the overriding method may leak the field (e.g., by assigning it to a static field).

Thus, *Cls* is augmented to include the classes that provide component functionality as well as all other classes transitively referenced. In the experiments presented in Section 5 we included all classes that were transitively referenced by *Cls*. This approach restricts analysis information to the currently “known world”—that is, the information may be invalidated in the future when new subclasses are added to *Cls*. Another approach is to change the analysis to make worst case assumptions for calls that may enter some unknown overriding methods. However, in this case, the analysis will be overly conservative and likely report fewer compositions. Thus, we believe that it is more useful to restrict the analysis to the known world; of course, the analysis user must be aware that the information is valid for the given set of known classes.

3. Points-to Analysis

Points-to analysis determines the set of objects that a given reference variable or a reference field may point to. This information has a wide variety of uses in software tools and optimizing compilers. In this paper, points-to information is used for ownership inference. It is needed to construct a graph that approximates all possible object graphs that can happen when arbitrary client code is built on top of *Cls*. There is a large body of work on points-to analysis with different trade-offs between cost and precision. In this paper we consider ownership inference based on the well-known Andersen-style flow- and context-insensitive points-to analysis for Java from [16].⁴

3.1. Points-to Analysis for Java

The points-to analysis is defined in terms of three sets. Set *R* is the set of locals, formals and static fields of reference type. Set *O* is the set of object names; the objects created at an allocation site s_i are represented by object name $o_i \in O$. Set *F* contains all instance fields in program

⁴ Flow-insensitive analyses do not take into account the flow of control between program points and are less precise and less expensive than flow-sensitive analyses. Context-sensitive analyses distinguish between different calling contexts of a method and are more precise and more expensive than context-insensitive ones.

classes. The analysis solution is a *points-to graph* where the edges represent the following "may-refer-to" relationships:

- Let $r \in R$ and $o \in O$. An edge $r \rightarrow o$ in the points-to graph means that at run time r may refer to some object that is represented by o .
- Let $f \in F$ be a reference instance field in objects represented by some $o \in O$. An edge $o \xrightarrow{f} o_2$ means that at run time field f of some object represented by o may refer to some object represented by o_2 .
- If o represents array objects, $o \Downarrow o_2$ shows that some element of some array represented by o may refer at run time to an object represented by o_2 .

The Andersen-style points-to analysis for Java from [16] is a relatively precise flow- and context-insensitive inclusion-based analysis. It propagates may-refer-to relationships by analyzing program statements. For example, when it analyzes statement " $p = q$ " it infers that p may refer to any object that q may refer to.

3.2. Fragment Points-to Analysis

Points-to analyses and Andersen's analysis in particular are typically designed as *whole-program analyses*; they take as input a complete program and produce points-to graphs that reflect relationships in the entire program. However, the problem considered in this paper requires points-to analysis of a partial program. The input is a set of classes Cls and the analysis needs to construct an approximate object graph that is valid across all possible executions of arbitrary client code built on top of Cls . To address this problem we make use of a general technique called *fragment analysis* due to Nasko Rountev [14, 17, 15]. Fragment analysis works on a program fragment rather than on a complete program; in our case the fragment is the set of classes Cls .

Initially, the fragment analysis produces an artificial `main` method that serves as a placeholder for client code written on top of Cls . Intuitively, the artificial `main` simulates the possible flow of objects between Cls and the client code. Subsequently, the fragment analysis attaches `main` to Cls and uses some whole-program analysis engine to compute a points-to graph which summarizes the possible effects of arbitrary client code. The fragment analysis approach can be used with a wide variety of points-to and class analyses; for the purposes of this paper we only consider fragment analysis used with the Andersen-style points-to analysis from [16].

The placeholder `main` method for the classes from Figure 3 is shown in Figure 4. The method contains variables for types from Cls that can be accessed by client code. The statements represent different possible interactions involving Cls ; their order is not relevant because the subsequent

```
void main() {
    ZipEntry ph_ZE;
    ZipInputStream ph_ZIS;
    ZipOutputStream ph_ZOS;
    ph_ZE = new ZipEntry();
    ph_ZIS = new ZipInputStream();
    ph_ZOS = new ZipOutputStream();
    ph_ZE.setCRC(0);
    ph_ZE = ph_ZIS.getNextEntry();
    ph_ZOS.putNextEntry(ph_ZE);
    ph_ZOS.closeEntry();
    ph_ZOS.finish(); }

```

Figure 4. Placeholder `main` method for `zip`.

whole-program analysis is flow-insensitive. Method `main` invokes all public methods from the classes in Cls designated as accessible.

The details of the fragment analysis will not be discussed here; they can be found in [17]. For the purposes of our analysis we discuss the *object reachability* [15] property of the results computed by the fragment analysis; this property is relevant for the analysis described in Section 4. Consider some client program built on top of Cls and an execution of this program (the program must satisfy the constraints discussed in Section 2.3). Let $r \in R$ be a variable declared in Cls and at some point during execution r is the start of a chain of object references that leads to some heap object. In the fragment analysis solution, there will be a chain of points-to edges that starts at r and leads to some object name o that represents the run-time object. A similar property holds if r is declared outside of Cls . In this case, in the fragment analysis solution, the starting point of the chain is the variable from `main` that has the same type as r .

We illustrate this property for our points-to analysis. Consider the example from Figures 3 and 4. There are three allocation sites in the `main` method; they are denoted by names `ZE1`, `ZIS1` and `ZOS1`. Name `byte[]` corresponds to the allocation site in class `InflaterInputStream`. There are three allocation sites in class `ZipInputStream`; they are denoted by names `CRC1`, `Inflater1` and `ZE2`. There are four allocation sites in class `ZipOutputStream`; they are denoted by `Vector1`, `Hashtable1`, `Deflater1` and `CRC2`. In addition, we consider the allocation sites in `Vector` (recall Figure 1), which are transitively reachable; they are denoted by `Object[]` and `VIter1`. The points-to graph computed by Andersen's analysis from the code in Figures 4, 3 and 1 is shown in Figure 5. Heap object names are underlined and reference variable names are prefixed by the name of their declaring method. For simplicity, implicit parameters `this` and object names `Inflater1`, `byte[]`, `Hashtable1` and `Deflater1` are not shown.

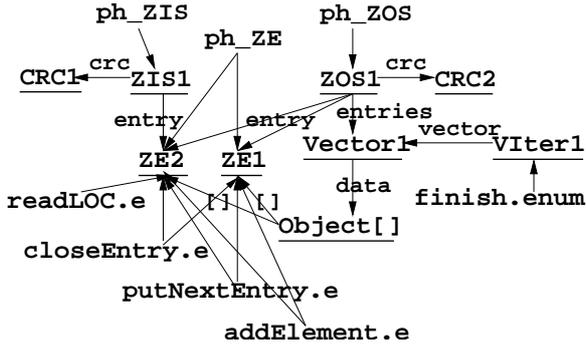


Figure 5. Points-to graph computed by the fragment points-to analysis.

4. Identifying Composition Relationships

We propose a novel analysis for ownership inference that is based on the output of the fragment points-to analysis. The ownership analysis constructs the *approximate object graph* Ag which approximates all possible run-time object graphs that can happen when client code is built on top of Cls . The analysis uses Ag to identify a *boundary* sub-graph rooted at o for each object name o ; the subgraph contains paths that are guaranteed to represent flow within the ownership boundary of o . Whenever the edge appears in the boundary of its source for *all* edges labeled with f , the relationship through f is identified as composition.

4.1. Approximate Object Graph

The nodes in Ag are taken from the set of object names O and the edges represent "may-access" relationships. Figure 6 outlines the construction of Ag given a points-to graph Pt . Set C_m denotes the set of object names that represent the contexts of invocation of method m . If m is an instance method or constructor, C_m is the points-to set of the implicit parameter `this` of m . If m is a static method C_m includes the union of the points-to sets of `this` for all instance methods or constructors that may call m (directly or through a sequence of static calls); it includes `root` if m is `main` or may be called from `main`.

Lines 1-2 account for edges due to flow from the contexts of the callee to the contexts of the caller. For example, at a constructor call new edges are added to Ag from each context enclosing the call to the name representing the newly created object. Similarly, at an instance call not through `this` new edges are added from each context enclosing the call to each returned object. Note that when the newly constructed object is immediately passed to another context (e.g., as in `new A(new B(...))`), or immediately returned to another context (e.g., as in `return new VIterator(this)`),

```

input  Stmt: set of statements  Pt:  $R \cup O \rightarrow \mathcal{P}(O)$ 
output  $Ag : O \rightarrow \mathcal{P}(O)$ 
[1] foreach
     $s : l = new C(\dots)$  s.t.  $l$  not immediately passed or
    immediately returned to another context,
     $s : l = r.m(\dots)$  s.t.  $r \neq this$ ,
     $s : l = r.f$  s.t.  $r \neq this$  and  $l$  assigned to a variable do
[2]  add  $\{c \rightarrow o_j \mid c \in C_{EnclMethod(s)} \wedge o_j \in Pt(l)\}$  to  $Ag$ 
    // add access edges due to flow from callees to callers
[3] foreach
     $s : l = new C(r)$ ,
     $s : l.m(r)$  s.t.  $l \neq this$ ,
     $s : l.f = r$  s.t.  $l \neq this$  do
[4]  add  $\{o_i \rightarrow o_j \mid o_i \in Pt(l) \wedge o_j \in Pt(r)\}$  to  $Ag$ 
    // add access edges due to flow from callers into callees
[5] foreach  $o_i \xrightarrow{f} o_j \in Pt$  do label with  $f$  each  $o_i \rightarrow o_j \in Ag$ 

```

Figure 6. Construction of Ag . $\mathcal{P}(X)$ denotes the power set of X . Ag is initially empty.

no new edges are added to that object from the contexts enclosing the constructor call. Also, at indirect read statements, no edges are added when variable l is not assigned or passed as an explicit or implicit argument later (e.g., it is used only to access instance or array fields such as in $x=l[i]$). This is consistent with the definition of the relaxed object graph in Section 2.1. Lines 3-4 account for edges due to flow from the contexts of the caller to the contexts of the callee. For example, at instance calls edges are added to each object in the points-to set of a reference argument, from each object in the points-to set of the receiver. Finally, line 5 labels the edges with the appropriate field identifier. For brevity, we omit discussion of static fields. The actual implementation creates edges from `root` to each object in the points-to set of a static field; the case is handled correctly by this algorithm and by the algorithm in Section 4.2.

We discuss the *reachability* property of the approximate object graph. Consider some client program built on top of Cls and an execution of this program (the program must satisfy the constraints discussed in Section 2.3). Let c be a context (i.e., `root` or a heap object) and at some point during execution c is the start of a chain in the relaxed object graph that leads to some heap object o^r . In Ag , there will be a chain of edges that starts at the representative of c and leads to the representative of o^r . Figure 7 shows the approximate object graph computed from the code on Figures 3, 4 and 1, and the points-to graph in Figure 5 (only object names from Figure 5 are shown). For the majority of edges inference is straight-forward. For example, edges $root \rightarrow ZIS1$, $root \rightarrow ZIS2$ and $root \rightarrow ZE1$ are due to the constructor calls in `main` and edges $ZIS1 \rightarrow CRC1$ and $ZIS1 \rightarrow ZE2$ are due to the constructor calls in

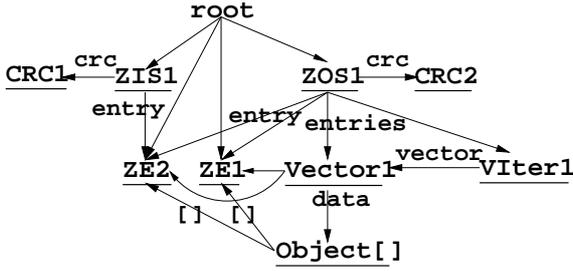


Figure 7. Approximate object graph computed by the algorithm in Figure 6.

```

input  Ag:  $O \rightarrow \mathcal{P}(O)$ 
output Bndry:  $O \rightarrow \mathcal{P}(O)$ 
initialize Bndry( $o_i$ )= $\{o_i\}$ 
         Forbid( $o_i$ )= $\{o_i \rightarrow o_j \text{ s.t. } \exists o_k \text{ s.t. } o_k \rightarrow o_i \wedge o_k \rightarrow o_j\}$ 
[1] while changes occur in Bndry
[2] foreach Bndry( $o_i$ ) and  $o \in Bndry(o_i)$ 
[3]  foreach  $o \rightarrow o_j$  not in Bndry( $o_i$ ) and not in Forbid( $o_i$ )
[4]   Wl= $\{\}$ , Tmp= $\{\}$ , inBndry=true
[5]   mark  $o \rightarrow o_j$ , and add to Wl and Tmp
[6]   while Wl not empty and inBoundary
[7]   remove  $o \rightarrow o_j$  from Wl
[8]   foreach  $o_k \rightarrow o_j \in Bndry(o_i)$  s.t.  $o_k \rightarrow o \wedge o_k \rightarrow o_j$ 
[9]    if  $o_k \rightarrow o_j$  unmarked, mark and add to Wl and Tmp
[10]  foreach  $o \rightarrow o_k$  s.t.  $o \rightarrow o_k \wedge o_k \rightarrow o_j$ 
[11]   if  $o \rightarrow o_k \in Bndry(o_i)$ 
[12]    if  $o_k \rightarrow o_j$  unmarked, mark and add to Wl and Tmp
[13]   else inBndry = false
[14] if inBndry
[15] add Tmp to Bndry( $o_i$ ), add Bndry( $o_j$ ) to Bndry( $o_i$ )

```

Figure 8. Ownership analysis.

class ZipInputStream. Edge $ZOS1 \rightarrow VIter1$ is due to call `enum=entries.elements()` in method `finish`. Edge $VIter1 \rightarrow Vector1$ is due to statement `return new VIterator(this)` in method `elements`; note that there is no edge $Vector1 \rightarrow VIter1$ due to this statement. Edge $root \rightarrow ZE2$ is due to statement `ph_ZE = ph_ZIS.getNextEntry()` in `main`, and edges $ZOS1 \rightarrow ZE2$ and $ZOS1 \rightarrow ZE1$ are due to statement `ph_ZOS.putNextEntry(ph_ZE)` in `main`. Edges $Object[] \rightarrow ZE2$ and $Object[] \rightarrow ZE1$ are due to flow at statement `data[at] = e` in `addElement`.

4.2. Ownership Boundary

The algorithm in Figure 8 takes Ag as input and outputs subgraphs $Bndry(o_i)$ for each object name. Subgraph

$Bndry(o_i)$ contains paths that are guaranteed to represent flow within the ownership boundary of an instance represented by o_i . More precisely, we have the following lemma. If a path $p: o_i \rightarrow o_1 \dots \rightarrow o_k \rightarrow o \in Bndry(o_i)$, then for each o_i^r (a run-time object represented by o_i) and o^r reachable from o_i^r on a path represented by p , we have (i) o_i^r dominates o^r and (ii) $Bndry(o_i)$ is closed with respect to o^r (i.e., the representative of every path from o_i^r to o^r is also in $Bndry(o_i)$). For example, the boundary of $ZOS1$ includes nodes $ZOS1, CRC2, Vector1, Object[]$ and $VIter1$ and the edges between them. It is easy to see that for every run-time path $ZOS1^r \rightarrow Vector1^r$, $ZOS1^r$ dominates $Vector1^r$. $Vector1^r$ is reachable along paths $ZOS1^r \rightarrow Vector1^r$ and $ZOS1^r \rightarrow VIter1^r \rightarrow Vector1^r$; the representatives of both paths are in $Bndry(ZOS1)$.

Below we briefly outline the algorithm and its correctness proof. The algorithm uses the fact that o_i^r flows from object o_i^r to some object o_k^r only if one of the following is true: (1) o_k^r has a handle to both o_i^r and o_j^r (and due to the reachability property Ag contains edges $o_k \rightarrow o_i, o_k \rightarrow o_j, o_i \rightarrow o_j$), or (2) o_i^r has a handle to both o_k^r and o_j^r (and Ag contains edges $o_i \rightarrow o_k, o_i \rightarrow o_j, o_k \rightarrow o_j$). This observation helps identify encapsulation more precisely. Suppose that our running example has another input stream object, created by `root` and denoted by name `ZIS2`. The relationship between `ZIS2` and its `crc` object would be represented by edge $ZIS2 \rightarrow CRC1$ in Figure 7. A naive algorithm may identify `root` as the dominator of the `crc` objects, and fail to identify the composition relationship between `ZipInputStream` and `CRC32`. In fact, the `CRC1` object is created and dominated by its enclosing `ZIS1` object because there is no o_k such that either o_k has handles to both `ZIS1` and `CRC1`, or `ZIS1` has handles to both o_k and `CRC1`; thus, the `CRC1` object created by the `ZIS1` object does not flow to or from any other context.

The algorithm builds the boundary of an object name o_i by adding edges. $Bndry(o_i)$ grows from zero to one edge, $o_i \rightarrow o_j$, when (i) there is no o_k that has handles to both o_i and o_j and (ii) there is no o_k such that o_i has handles to both o_k and o_j , and o_k has a handle to o_j . The first condition is guaranteed by the check at line 3 and the second condition is guaranteed by the fact that the flag variable `inBndry` stays true only if the loop at lines 10-13 is skipped. Thus, an edge $o_i \rightarrow o_j$ is added to the empty boundary of o_i only when it is guaranteed that the o_i object accesses the o_j object exclusively (i.e., no other object has a handle to it). Examples of such edges are $ZIS1 \rightarrow CRC1$ and $ZOS1 \rightarrow CRC2$. Clearly, the conditions of the lemma hold in this case.

Consider an edge $o \rightarrow o_j$ (examined at lines 3-15). Consider some client program built on top of `Cls` and an execution of this program (the program must satisfy the constraints discussed in Section 2.3). Let o_i^r be any run-time object represented by o_i and o^r be an object reachable on

a path $o_i^r \rightarrow \dots o_k^r \rightarrow o^r$ whose representative is in the boundary of o_i . We may assume that (i) o_i^r dominates o^r and (ii) $Bndry(o_i)$ is closed w.r.t. o^r . Let o^r refer to an instance o_j^r . We need to examine all o_k^r such that o_j^r may flow to or from o_k^r (i.e., there is an edge $o_k^r \rightarrow o_j^r$). If all these o_k^r are dominated by o_i^r then o_i^r dominates o_j^r . If the boundary is closed w.r.t. each o_k^r adding edges $o \rightarrow o_j$ and $o_k \rightarrow o_j$ ensures that the boundary is closed w.r.t. o_j^r .

At lines 4-5 the algorithm marks the edge as visited in the current iteration, sets the *inBndry* flag to true, and initializes the worklist and the closure set *Tmp* (discussed later) to $o \rightarrow o_j$. Object o_j^r flows from o^r into some o_k^r when one of the following conditions is true. First, o_k^r has handles to both o^r and o_j^r (e.g., o_j^r may be returned from a method invoked on o^r). Since all paths to $o_i^r \rightarrow \dots \rightarrow o_k^r$ are subpaths of paths to o^r , their representatives must be in $Bndry(o_i)$; thus, o_k^r is (i) dominated by o_i^r and (ii) $Bndry(o_i)$ is closed w.r.t. o_k^r . This case is examined at lines 8-9 and $o_k \rightarrow o_j$ is added to the worklist; it is examined in a subsequent iteration of the while loop in order to find the representatives of the objects that o_j^r may flow to from o_k^r . In addition, $o_k \rightarrow o_j$ is added to *Tmp*, the closure set of $o \rightarrow o_j$ —adding the edges in the closure set to the boundary (line 15) ensures that the boundary is closed w.r.t. each o_j^r . Second, o_j^r may flow from o^r into some o_k^r such that o^r has handles to both o_k^r and o_j^r . Since $o \rightarrow o_k \in Bndry(o_i)$ we have that (i) o_k^r is dominated by o_i^r and (ii) the boundary is closed w.r.t. o_k^r . This case is examined at the loop at lines 10-13 and appropriate $o_k \rightarrow o_j$ are added to the worklist and to the closure set. Finally, if the edges in the closure set form appropriate paths, *inBndry* is true and at line 15 the algorithm adds the closure set *Tmp* and the boundary of o_j to the boundary of o_i .

We briefly illustrate the algorithm on our running example. Consider edge $ZIS1 \rightarrow CRC1$. The algorithm skips the loop on lines 8-9 because set $Bndry(ZIS1)$ is empty; it skips the loop on lines 10-13 as well because there is no o_k with handles to both $ZIS1$ and $CRC1$. At line 15, edge $ZIS1 \rightarrow CRC1$ is added to $Bndry(ZIS1)$. Edges $ZOS1 \rightarrow CRC2$, $ZOS1 \rightarrow VIter1$ and $Vector1 \rightarrow Object[]$ are processed analogously. For edge $ZOS1 \rightarrow Vector1$ the algorithm goes through lines 10-13 with $o \rightarrow o_k$ being $ZOS1 \rightarrow VIter1$. It determines that $ZOS1 \rightarrow VIter1$ is already in the boundary of $ZOS1$; thus, the vector object can be accessed only by its enclosing output stream object and the iterator object which is enclosed in that output stream object. Thus we have the following boundary graphs: $Bndry(ZIS1) = \{ZIS1 \rightarrow CRC1\}$, $Bndry(Vector1) = \{Vector1 \rightarrow Object[]\}$ and $Bndry(ZOS1) = \{ZOS1 \rightarrow CRC2, ZOS1 \rightarrow Vector1, ZOS1 \rightarrow VIter1, Vector1 \rightarrow Object[], VIter1 \rightarrow Vector1\}$.

A corollary of the lemma is that whenever we have an

edge $o_i \rightarrow o_j \in Bndry(o_i)$ each o_i^r owns the o_j^r instances that it may refer to. If for every edge labeled with f we have $o \xrightarrow{f} o' \in Bndry(o)$ the analysis identifies one-to-one implementation-level composition or collection ownership.

4.3. Composition Relationships

Let f be an instance field of reference type in *Cls*. Consider all field edges $o \xrightarrow{f} o'$ in *Pt*. Clearly, if for each edge $o \rightarrow o' \in Bndry(o)$, the ownership property holds for f and the analysis identifies one-to-one implementation-level composition or collection ownership.

Consider an edge $o \xrightarrow{f} o'$ where f is a field of collection type. The set of instances stored in collections represented by o' is approximated by set $Stored(o')$, the union of the points-to sets of the variables passed as appropriate actual arguments to standard put methods invoked on o' . For example, for Figures 3 and 4, set $Stored(Vector1)$ is the points-to set of variable *e* in `putNextEntry` which is the actual argument to the standard put method of `Vector`, `addElement`; therefore $Stored(Vector1) = \{ZE1, ZE2\}$. Similarly, set $Stored(Hashtable1)$ is the points-to set of *e* in `putNextEntry`; therefore we have $Stored(Hashtable1) = \{ZE1, ZE2\}$. We extend the notation to have $Stored(f)$ denote the union of the sets $Stored(o')$ where o' is a collection object stored in f . If for each path p through an edge $o \rightarrow o'$ to a name in $Stored(f)$ we have $p \in Bndry(o)$, the analysis identifies a one-to-many composition relationship through f ; otherwise the analysis identifies a one-to-many association relationship.

It remains to identify the least common supertype C of the objects stored in collections fields as the analysis needs to infer an association between the enclosing class of f and C . As explained in Section 2.2 our analysis uses a new type system to infer the least common supertype of all types in $Stored(f)$, then maps the new type to a Java type. If a Java type does not exist, the analysis uses `Object`, thus losing precision. In our example, the least common supertype of all instances stored in `Vector1` is easily identified to be `ZipEntry`; thus, the analysis infers that there is a regular one-to-many association between `ZipOutputStream` and `ZipEntry` through field entries.

5. Experimental Study

The goal of the study is to address two questions. First, how often does our analysis discover implementation-level composition? Second, how *imprecise* the analysis is—that is, how often it misses implementation-level composition?

For the experiments we used several Java components from the standard library packages `java.text` and

`java.util.zip` [15]. The components are described briefly in the first two columns of Table 1. Each component contains the set of classes in *Cls* (these are the classes that provide certain functionality plus all other classes that are directly or transitively referenced by them); the number of classes in each component is shown in column (3). We considered all reference instance fields in the classes in *Cls* that provide the component functionality; this number is given in column (4).

5.1. Results

We applied the algorithm described earlier in order to determine which fields accounted for composition relationships. The results are given in the last two columns of Table 1. Column (5) shows how many of the fields from column (4) are identified as one-to-one compositions and column (3) shows how many of the fields are identified as owned collections (i.e., arrays and standard `java.util` collections).

On average, the analysis reported 30% one-to-one compositions and 10% owned collections—that is, 40% of the reference instance fields account for representation that is not being exposed outside of its enclosing object. The owned collections in component `zip` are analogous to the `Vector` and `Hashtable` in our running example; they store exposed `ZipEntries` and account for one-to-many association relationships. Two of the owned collections for component `collator` and one of the owned collections for `date` accounted for one-to-many compositions. The remaining owned collections are arrays of simple type.

5.2. Analysis Precision

The issue of analysis precision is of crucial importance for software tools. If an analysis is imprecise, it may report that the relationship between two classes is not a composition while in reality it is, or that a collection is not owned while in reality it is owned (i.e., the analysis reports that certain representation may be exposed while in fact it is not). Such information is not useful and may confuse the user and even render the tool unusable. For example, if a user attempts to ensure the consistency between the code and the composition relationships in UML design class diagrams, imprecision will mean that a large chunk of code will have to be examined manually. Since imprecision results in waste of human time, analysis designers must carefully and precisely identify and evaluate any sources of imprecision.

In our experiments, we examined the fields in *Cls* that were not identified as compositions or owned collections. We attempted to prove that it was possible to write client code for which some object stored in such a field would be exposed (i.e., it would not be owned by its enclosing ob-

ject in accordance with the ownership model in Section 2.1). In all cases, except one, we were able to prove that the object was exposed. Thus, the analysis achieves almost perfect precision.

Field `defaultCenturyStart` in component `date` accounted for the one case of imprecision. The imprecision was due to context-insensitive object naming in the points-to and composition analyses. The object stored in the field comes from a call to a method `getTime` which creates and immediately returns a `Date` object. Although the `Date` object stored in `defaultCenturyStart` does not flow out of its enclosing `SimpleDateFormat` object, other `Date` objects created by `getTime` in `SimpleDateFormat` are being returned (i.e., there are edges in *Ag* to the `SimpleDateFormat` object and the only representative of `Date`). This imprecision may be resolved by using an analysis that employs more precise object naming. In the case of `getTime` it may distinguish the `Date` objects for different call sites of `getTime`; the target of the `defaultCenturyStart` edge would be a separate `Date` object that does not flow out and the ownership inference algorithm will correctly identify that there is a composition relationship through this field. However, it is unclear whether a more precise context-sensitive points-to analysis will result in substantial benefits for the ownership and composition analyses.

5.3. Conclusions

Our results indicate that the ownership model captures conceptual composition relationships appropriately—we encountered several cases when values of private fields were stored in other parts of the object representation. Thus, a simpler model based on exclusive ownership—that is, a model which requires that an owned object is in exclusive relationship with its owner, would not have been sufficient to identify compositions. The results also show that composition relationships occurs often. Therefore, the analysis can provide useful information for reverse engineering tools. It is important that highly precise information can be obtained with practical analysis—the combined running time of the points-to and composition inference analyses does not exceed 10 seconds on any component (based on the average of three runs on a Sun Fire 380R). Of course, these results need to be reconfirmed on more components.

6. Related Work

Recent work by Guéhéneuc and Albin-Amiot [9] presents definitions and identification algorithms for implementation-level association, composition and aggregation relationships. Our work focuses more closely on compositions and differs substantially from [9] in

(1)Component	(2)Functionality	(3)#Classes	(4)#Fields	Compositions			
				(5)#One-to-one		(6)#Owned collections	
				Analysis	Perfect	Analysis	Perfect
gzip	GZIP IO streams	199	7	4(57%)	4(57%)	0(0%)	0(0%)
zip	ZIP IO streams	194	10	3(30%)	3(30%)	2(20%)	2(20%)
checked	IO streams with checksums	189	2	0(0%)	0(0%)	0(0%)	0(0%)
collator	text collation	203	24	10(42%)	10(42%)	6(25%)	6(25%)
date	date formatting	205	20	3(15%)	4(20%)	5(25%)	5(25%)
number	number formatting	198	3	2(67%)	2(67%)	0(0%)	0(0%)
boundary	iteration over boundaries in text	199	7	0(0%)	0(0%)	0(0%)	0(0%)
Average				30%	31%	10%	10%

Table 1. Java components and implementation-level compositions.

both the definition of implementation-level composition and in the proposed identification algorithm. The definition of composition in [9] is based on exclusive ownership. This may not be sufficient to model such commonly used object-oriented patterns and idioms such as iterators, composites, decorators, and factories [7] as well as the common situation when instance fields refer to owned objects that are stored in owned collections or temporarily accessed by other parts of the representation of the owner. Our definition is based on the owners-as-dominators model which does not require exclusive relationship with the owner; as observed by us and other researchers [4, 13], this model captures well the notion of composition in modeling [18].

Most importantly, we present an identification algorithm that may be more appropriate. Guéhéneuc and Albin-Amiot propose the use of dynamic analysis, but point out serious disadvantages. First, dynamic analysis is slow, second, it requires a complete program, and third, the results that are obtained may be incomplete because they are based on particular runs of particular clients of the component. Our detection algorithm is based on practical static analysis that works on incomplete programs and produces a solution that is valid over all unknown clients of the component.

Work in [10] and [19] addresses the issue of recovering one-to-many associations through containers, since reverse engineering tools typically lose the association between the enclosing class and the class whose instances are stored in the container field (recall the `entries` field of `Vector` type in Figure 3). Identification of composition is not addressed in these papers. Although our work focuses on identification of composition, our algorithm identifies one-to-many associations as well, and the algorithm is more general than those in [10] and [19].

Ownership type systems disallow certain accesses of object representation [12, 4, 3, 1, 2]. These systems require type annotations and typically do not include automatic inference algorithms or empirical investigations. In contrast,

we infer ownership automatically and present an empirical study of the effectiveness of our approach; we believe that our analysis can be usefully incorporated in software tools for reverse engineering of class diagrams from Java code. The only type annotation inference analysis that we are aware of is given by Aldridge et al. [1] for the purposes of alias understanding. Similarly to the work by Guéhéneuc and Albin-Amiot [9], the owned annotation is used only when the analysis is able to prove *exclusive* ownership; in the majority of cases it infers alias parameters. Our work focuses on a different problem, composition inference, and infers ownership using a more general ownership model that captures better the notion of composition in modeling. In addition, our algorithm may scale better.

7. Conclusions and Future Work

This work presents an approach for performing analysis that identifies composition relationships in Java components. We define an ownership-based implementation-level composition model and a static analysis that infers composition relationships in incomplete programs. Our empirical study indicates that (i) the ownership-based model captures well the notion of composition in modeling and (ii) implementation-level compositions occur often and almost *all* such compositions can be identified. Clearly, no definitive conclusions can be drawn from these limited experiments. In our future work we plan to focus on further empirical investigation.

References

- [1] J. Aldrich, V. Kostadinov, and C. Chambers. Alias annotations for program understanding. In *OOPSLA*, 2002.
- [2] C. Boyapati, B. Liskov, and L. Shriram. Ownership types for object encapsulation. In *POPL*, 2003.
- [3] D. Clarke and S. Drossopoulou. Ownership, encapsulation and the disjointness of type and effect. In *OOPSLA*, 2002.

- [4] D. Clarke, J. Potter, and J. Noble. Ownership types for flexible alias protection. In *OOPSLA*, 1998.
- [5] A. Donovan, A. Kiezun, M. S. Tschantz, and M. D. Ernst. Converting Java programs to use generic libraries. In *OOPSLA*, 2004.
- [6] M. Fowler. *UML Distilled Third Edition*. Addison-Wesley, 2004.
- [7] E. Gamma, R. Helm, R. Johnson, and J. Vlissides. *Design Patterns: Elements of Reusable Object-Oriented Software*. Addison-Wesley, 1995.
- [8] J. Gosling, B. Joy, G. Steele, and G. Bracha. *The Java Language Specification*. Addison-Wesley, 2nd edition, 2000.
- [9] Y.-G. Guéhéneuc and H. Albin-Amiot. Recovering binary class relationships: Putting icing on the UML cake. In *OOPSLA*, 2004.
- [10] D. Jackson and A. Waingold. Lightweight recovery of object models from bytecode. In *ICSE*, 1999.
- [11] C. Larman. *Applying UML and Patterns*. Prentice Hall, 2nd edition, 2002.
- [12] J. Noble, J. Vitek, and J. Potter. Flexible alias protection. In *ECOOP*, 1998.
- [13] J. Potter, J. Noble, and D. Clarke. The ins and outs of objects. In *Australian Software Engineering Conference*, 1998.
- [14] A. Rountev. *Dataflow Analysis of Software Fragments*. PhD thesis, Rutgers University, 2002.
- [15] A. Rountev. Precise identification of side-effect free methods. In *ICSM*, 2004.
- [16] A. Rountev, A. Milanova, and B. G. Ryder. Points-to analysis for Java using annotated constraints. In *OOPSLA*, 2001.
- [17] A. Rountev, A. Milanova, and B. G. Ryder. Fragment class analysis for testing of polymorphism in Java software. *IEEE TSE*, 30(6):372–386, June 2004.
- [18] J. Rumbaugh, I. Jacobson, and G. Booch. *The Unified Modeling Language Reference Manual*. Addison-Wesley, 2nd edition, 2004.
- [19] P. Tonella and A. Potrich. Reverse engineering of the uml class diagram from c++ code in presence of weakly typed containers. In *ICSM*, pages 376–385, 2001.
- [20] D. von Dincklage and A. Diwan. Converting Java classes to use generics. In *OOPSLA*, 2004.