

Reengineering Class Hierarchies Using Concept Analysis

GREGOR SNELTING

Universität Passau

FRANK TIP

IBM T.J. Watson Research Center

Abstract

A new method is presented for analyzing and reengineering class hierarchies. In our approach, a class hierarchy is processed along with a set of applications that use it, and a fine-grained analysis of the access and subtype relationships between objects, variables and class members is performed. The result of this analysis is again a class hierarchy, which is guaranteed to be behaviorally equivalent to the original hierarchy, but in which each object only contains the members that are required. Our method is semantically well-founded in *concept analysis*: the new class hierarchy is a minimal and maximally factorized *concept lattice* that reflects the access and subtype relationships between variables, objects and class members.

The method is primarily intended as a tool for finding imperfections in the design of class hierarchies, and can be used as the basis for tools that largely automate the process of reengineering such hierarchies. The method can also be used as a space-optimizing source-to-source transformation that removes redundant fields from objects.

A prototype implementation for Java has been constructed, and used to conduct several case studies. Our results demonstrate that the method can provide valuable insights into the usage of the class hierarchy in a specific context, and lead to useful restructuring proposals.

1 Introduction

Designing a class hierarchy is hard, because it is not always possible to anticipate how a hierarchy will be used by an application. This is especially the case when a class hierarchy

⁰A preliminary version of parts of this article appeared in the Proceedings of the ACM SIGSOFT Symposium on the Foundations of Software Engineering, 1998 [30].

Authors' addresses: Gregor Snelting, Fakultät für Mathematik und Informatik, Universität Passau, Innstr. 33, 94032 Passau, Germany; Frank Tip, IBM T.J. Watson Research Center, P.O. Box 704, Yorktown Heights, NY 10598, USA

```

class String { /* details omitted */ };
class Address { /* details omitted */ };
enum Faculty { Mathematics, ComputerScience };
class Professor; /* forward declaration */

class Person {
public:
    String name;
    Address address;
    long socialSecurityNumber;
};

class Student : public Person {
public:
    Student(String sn, Address sa, int si){
        name = sn; address = sa; studentId = si;
    };
    void setAdvisor(Professor *p){
        advisor = p;
    };
    long studentId;
    Professor *advisor;
};

class Professor : public Person {
public:
    Professor(String n, Faculty f, Address wa){
        name = n; faculty = f;
        workAddress = wa;
        assistant = 0; /* default: no assistant */
    };
    void hireAssistant (Student *s){
        assistant = s;
    };
    Faculty faculty;
    Address workAddress;
    Student *assistant; /* either 0 or 1 assistants */
};

```

(a)

```

int main(){
    String s1name, p1name;
    Address s1addr, p1addr;
    Student* s1 = /* Student1 */
        new Student(s1name,s1addr,12345678);
    Professor *p1 = /* Professor1 */
        new Professor(p1name,Mathematics,p1addr);
    s1->setAdvisor(p1);
    return 0;
}

```

(b)

```

int main(){
    String s2name, p2name;
    Address s2addr, p2addr;
    Student* s2 = /* Student2 */
        new Student(s2name,s2addr,87654321);
    Professor *p2 = /* Professor2 */
        new Professor(p2name, ComputerScience, p2addr);
    p2->hireAssistant(s2);
    return 0;
}

```

(c)

Figure 1: Example: relationships between students and professors. (a) Class hierarchy for expressing associations between students and professors. (b) Example program using the class hierarchy of Figure 1(a). (c) Another example program using the class hierarchy of Figure 1(a).

is developed as a library, and designed independently from the applications that use it. Ongoing maintenance, in particular ad-hoc extensions of the hierarchy, will further increase the system’s entropy. As typical examples of inconsistencies that may arise, one might think of:

- A class C may contain a member m not accessed in any C -instance, an indication that m may be removed, or moved into a derived class.
- Different instances of a given class C may access different subsets of C ’s members, an indication that it might be appropriate to split C into multiple classes.

In this paper, we present a method for analyzing the usage of a class hierarchy based on *concept analysis* [37]. Our approach comprises the following steps. First, a table is constructed that precisely reflects the usage of a class hierarchy. In particular, the table makes explicit relationships between the types of variables and class members such as “the type of x must be a base class of the type of y ”, and “member m must occur in a base class of the type of variable x ” are encoded in the table. From the table, a *concept lattice* is derived, which factors out information that variables or members have in common. We will show how the concept lattice can provide valuable insight into the design of the class hierarchy, and how it can serve as a basis for automated or interactive restructuring tools for class hierarchies. The examples presented in this paper are written in C++ or Java, but our approach is applicable to other object-oriented languages as well.

Our method can analyze a class hierarchy along with any number of programs that use it, and provide the user with either a combined view reflecting the usage of the hierarchy by the entire set of programs, or with individual views that clarify how each application uses the hierarchy. Analyzing a class hierarchy without any accompanying applications (such as a class library) is also possible, and can be useful to study the internal dependences inside class definitions.

1.1 A motivating example

Consider the example of Figure 1, which is concerned with relationships between students and professors. Figure 1(a) shows a class hierarchy, in which a class `Person` is defined that contains a person’s `name`, `address`, and `socialSecurityNumber`. Classes `Student` and `Professor` are derived from `Person`. `Students` have an identification number (`studentId`), and a thesis `advisor` if they are graduate students. A constructor is provided for initializing `Students`, and a method `setAdvisor` for designating a `Professor` as an advisor. `Professors` have a `faculty` and a `workAddress`, and a professor may hire a student as a teaching `assistant`. A constructor is provided for initialization, and a method `hireAssistant` for hiring a `Student` as an assistant. Details for classes `Address` and `String` are not provided; in the subsequent analysis these classes will be treated as “atomic” types and we will not attempt to analyze them.

Figure 1(b) and (c) show two programs that use the class hierarchy of Figure 1(a). In the first program, a student and a professor are created, and the professor is made the

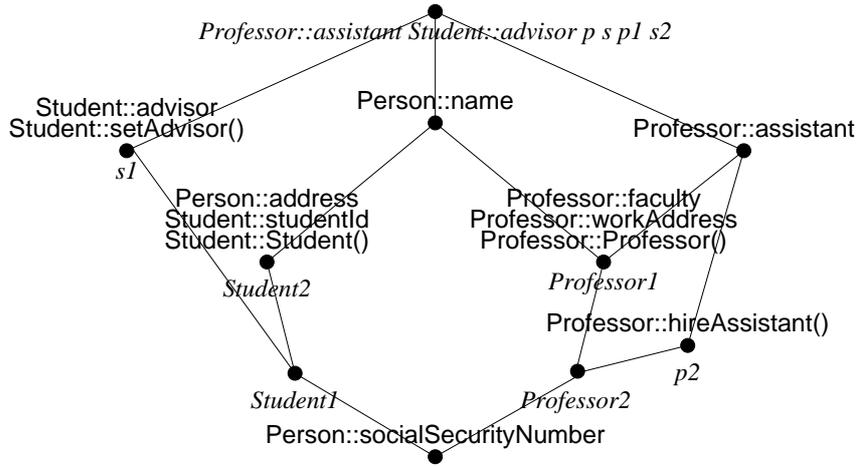


Figure 2: Lattice for Student/Professor example.

student’s advisor. The second program creates another student and professor, and here the student is made the professor’s assistant. The example is certainly not perfect C++ code, but looks reasonable enough at first glance.

Figure 2 shows the lattice computed by our method for the class hierarchy and the two example programs of Figure 1. Ignoring a number of details, the lattice may be interpreted as follows:

- The lattice elements (concepts) may be viewed as *classes* of a restructured class hierarchy that precisely reflects the usage of the original class hierarchy by the client programs.
- The ordering between lattice elements may be viewed as *inheritance* relationships in the restructured class hierarchy.
- A variable v has type C in the restructured class hierarchy if v occurs immediately *below* concept C in the lattice.
- A member m occurs in class C if m appears *directly above* concept C in the lattice.

Examining the lattice of Figure 2 according to this interpretation reveals the following interesting facts¹:

- Data member `Person::socialSecurityNumber` is never accessed, because no variable appears below it. This illustrates situations where subclassing is used to inherit the functionality of a class, but where some of that functionality is not used.

¹The labels `Student1`, `Professor1`, `Student2`, and `Professor2` that appear in the lattice represent the types of the heap objects created by the example programs at various program points (indicated in Figures 1(b) and (c) using comments).

- Data member `Person::address` is only used by students, and not by professors (for professors, the data member `Professor::workAddress` is used instead, perhaps because their home address is confidential information). This illustrates a situation where the member of a base class is used in some, but not all derived classes.
- No members are accessed from parameters `s` and `p`, and from data members `advisor` and `assistant`. This is due to the fact that no operations are performed on a student's advisor, or on a professor's assistant. Such situations are typical of redundant, incomplete, or erroneous code and should be examined closely.
- The analyzed programs create professors who hire assistants (`Professor2`), and professors who do not hire assistants (`Professor1`). This can be seen from the fact that method `Professor::hireAssistant()` appears above the concept labeled `Professor2`, but not above the concept labeled `Professor1`.
- There are students with advisors (`Student1`) and students without advisors (`Student2`). This can be seen from the fact that `Student::setAdvisor` appears above the concept labeled `Student1`, but not above the concept labeled `Student2`.
- Class `Student`'s constructor does not initialize the `advisor` data member. This can be seen from the fact that data member `Student::advisor` does not appear above method `Student::Student()` in the lattice².

One can easily imagine how the above information might be used as the basis for restructuring the class hierarchy. One possibility would be for a tool to automatically generate restructured source code from the information provided by the lattice, similar to the approach taken in [35, 36]. However, from a redesign perspective, we believe that an interactive approach would be more appropriate. For example, the programmer doing the restructuring job may decide that the data member `socialSecurityNumber` should be retained in the class hierarchy because it may be needed later. In the interactive tool we envision, one could indicate this by *moving up* in the lattice the attribute under consideration, `socialSecurityNumber`. The reengineer may also decide that certain fine distinctions in the lattice are unnecessary. For example, one may decide that it is not necessary to distinguish between professors that hire assistants, and professors that don't. In an interactive tool, this distinction could be removed by *merging* the concepts for `Professor1` and `Professor2`.

Another useful capability of an interactive tool would be to associate names with lattice elements. When the programmer is done manipulating the lattice, these names could be used as class names in the restructured hierarchy when the restructured source code is generated. For example, using the information provided by the lattice, the programmer may determine that `Student` objects on which the `setAdvisor` method is invoked are graduate students, whereas `Student` objects on which this method is not called are undergraduates. Consequently, he may decide to associate the names `Student` and `GraduateStudent` with the concepts labeled `Student2` and `Student1`, respectively.

²`Student::Student()` also represents the `this`-pointer of the method.

1.2 Organization of this paper

The remainder of this paper is organized as follows. Section 2 briefly reviews the relevant parts of the theory of concept analysis. In Section 3 we define the objects and attributes in our domain, which correspond to the rows and columns of the tables. The process of constructing tables is presented in Section 4, while Section 5 discusses important properties of the lattice, in particular behavioral equivalence. Section 6 presents extensions for constructs such as type casts. In Section 7, we discuss how the information provided by the lattice can reveal problems in the design of class hierarchies, and how the lattice can be used as a basis for interactive restructuring tools. Section 8 describes our prototype implementation for Java in some detail. Section 9 discusses several case studies. Section 10 discusses related work. Finally, conclusions and directions for future work are presented in Section 11.

2 Concept analysis

Concept analysis provides a way to identify groupings of *objects* that have common *attributes*. The mathematical foundation was laid by Birkhoff in 1940 [5]. Birkhoff proved that for every binary relation between certain objects and attributes, a lattice can be constructed that provides remarkable insight into the structure of the original relation. The lattice can always be transformed back to the original relation, hence concept analysis is similar in spirit to Fourier analysis.

Later, Wille and Ganter elaborated Birkhoff's result and transformed it into a data analysis method [37, 13]. Since then, it has found a variety of applications, including analysis of software structures [16, 28, 18, 27, 14, 29, 15].

2.1 Relations and their lattices

Concept analysis starts with a relation, or boolean table, T between a set of *objects* \mathcal{O} and a set of *attributes* \mathcal{A} , hence $T \subseteq \mathcal{O} \times \mathcal{A}$.

For any set of objects $O \subseteq \mathcal{O}$, their set of common attributes is defined as

$$\sigma(O) = \{a \in \mathcal{A} \mid \forall o \in O : (o, a) \in T\}$$

For any set of attributes $A \subseteq \mathcal{A}$, their set of common objects is

$$\tau(A) = \{o \in \mathcal{O} \mid \forall a \in A : (o, a) \in T\}$$

A pair (O, A) is called a *concept* if

$$A = \sigma(O) \text{ and } O = \tau(A)$$

Informally, such a concept corresponds to a *maximal rectangle* in the table T : any $o \in O$ has all attributes in A , and all attributes $a \in A$ fit to all objects in O . It is important to

	small	medium	large	near	far	moon	no moon
Mercury	×			×			×
Venus	×			×			×
Earth	×			×		×	
Mars	×			×		×	
Jupiter			×		×	×	
Saturn			×		×	×	
Uranus		×			×	×	
Neptune		×			×	×	
Pluto	×				×	×	

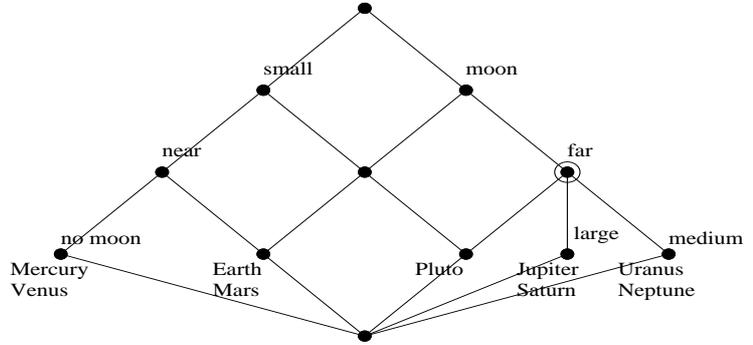


Figure 3: Example table and associated concept lattice.

note that concepts are invariant against row or column permutations in the table. The set of all concepts of a given table forms a partial order via

$$(O_1, A_1) \leq (O_2, A_2) \iff O_1 \subseteq O_2 \iff A_1 \supseteq A_2$$

Birkhoff proved that the set of concepts constitutes a complete lattice, the *concept lattice* $\mathcal{L}(T)$. For two elements (O_1, A_1) and (O_2, A_2) in the concept lattice, their infimum or *meet* is defined as

$$(O_1, A_1) \wedge (O_2, A_2) = (O_1 \cap O_2, \sigma(O_1 \cap O_2))$$

and their supremum or *join* as

$$(O_1, A_1) \vee (O_2, A_2) = (\tau(A_1 \cap A_2), A_1 \cap A_2)$$

A concept $c = (O, A)$ has *extent* $\text{ext}(c) = O$ and *intent* $\text{int}(c) = A$. In our figures, a lattice element (concept) c is labeled with attribute $a \in \mathcal{A}$, if it is the *largest* concept with a in its intent, and it is labeled with an object $o \in \mathcal{O}$, if it is the *smallest* concept with o in its extent. The (unique) lattice element labeled with a is denoted $\mu(a)$, and the (unique) lattice element labeled with o is denoted $\gamma(o)$. Thus

$$\mu(a) = \bigvee \{c \in \mathcal{L}(T) \mid a \in \text{int}(c)\}, \quad \gamma(o) = \bigwedge \{c \in \mathcal{L}(T) \mid o \in \text{ext}(c)\}$$

The following fundamental property establishes the connection between a table and its lattice, and shows that they can be reconstructed from each other:

$$(o, a) \in T \iff \gamma(o) \leq \mu(a)$$

Hence, the attributes of object o are those which appear *above* o , and all objects that appear *below* a have attribute a . Consequently, join points (suprema) in the lattice indicate that

certain objects have attributes in common, while meet points (infima) show that certain attributes fit to common objects. In other words, join points factor out common attributes, while meet points factor out common objects. Thus, the lattice uncovers a hierarchy of conceptual clusters that was implicit in the original table.

Figure 3 shows a table and its lattice (taken from [10]). The element labeled *far* corresponds to the maximal rectangle indicated in the table. This element is the supremum of all elements with *far* in their intent: *Pluto*, *Jupiter*, *Saturn*, *Uranus*, *Neptune* are below *far* in the lattice, and the table confirms that these (and no other) planets are indeed far away.

2.2 Implications

A table and its lattice are alternate views on the same information, serving different purposes and providing different insights. There is yet another view: a set of *implications*. Let $A, B \subseteq \mathcal{A}$ be two sets of attributes. We say that A *implies* B , iff any object with the attributes in A also has the attributes in B :

$$A \rightarrow B \iff \forall o \in \mathcal{O} : (\forall a \in A : (o, a) \in T) \Rightarrow (\forall b \in B : (o, b) \in T)$$

For $B = \{b_1, \dots, b_k\}$, $A \rightarrow B$ holds iff $A \rightarrow b_i$ for all $b_i \in B$.³ Implications show up in the lattice as follows: $A \rightarrow b$ holds iff $\bigwedge \{\mu(a) \mid a \in A\} \leq \mu(b)$. Informally, implications between attributes can be found along upward paths in the lattice. In the example of Figure 3, we have that $\mu(\text{far}) \leq \mu(\text{moon})$, which can be read as $\text{far} \rightarrow \text{moon}$, or “A planet which is far away has a moon”. Other examples of implication are $\text{nomoon} \rightarrow \text{near}, \text{small}$; or $\text{near}, \text{far} \rightarrow \text{large}$ (the latter implication being true because its premise is contradictory).

There is a minimal set of implications, from which all other valid implications can be derived: the *implication base* $\mathcal{I}(T)$. For the example, it consists of 10 implications, including $\text{far} \rightarrow \text{moon}$ and $\text{no moon} \rightarrow \text{near}, \text{small}$. Non-base implications such as $\text{far}, \text{small} \rightarrow \text{moon}, \text{small}$ or $\text{no moon} \rightarrow \text{near}$ can be derived by propositional logic.

Often, some implications are known to hold *a priori*. Such background knowledge can easily be integrated into a given table. An implication $x \rightarrow y$ can be enforced by copying the entries from the x column to the y column, and will cause $\mu(x) \leq \mu(y)$ in $\mathcal{L}(T)$. A general implication $A \rightarrow B$ can be enforced by copying the intersection of the A columns to all B columns.

2.3 Lattice construction

The table T , the lattice $\mathcal{L}(T)$, and the implication basis $\mathcal{I}(T)$ represent very different views onto the same information, but can be transformed into each other⁴; furthermore, background knowledge, given as a set of implications, may be added. In this section, we

³ We will usually write $a_1, \dots, a_n \rightarrow b_1, \dots, b_m$ instead of $\{a_1, \dots, a_n\} \rightarrow \{b_1, \dots, b_m\}$.

⁴Note the analogy to Fourier analysis: a function and its Fourier transform are very different representations of the same information, but can be transformed into each other.

will present a short description of the most important transformation: the computation of the concept lattice for a given table.

Ganter’s algorithm for lattice construction utilizes the fact that $C = \sigma \circ \tau$ (as well as $C' = \tau \circ \sigma$) is a closure operator on $2^{\mathcal{O}}$: it is extensive ($O \subseteq C(O)$), idempotent ($C(C(O)) = C(O)$), and monotone ($O \subseteq O' \Rightarrow C(O) \subseteq C(O')$). $C(O)$ determines the largest object set with the same common attributes as O . It turns out that the lattice elements’ extents are precisely the closed sets under C . If we have computed all the extents (that is, computed the closure system $\{C(O) \mid O \subseteq \mathcal{O}\}$), the corresponding intents are determined using σ , and the lattice, together with its partial order as defined above, is complete.

Ganter’s algorithm requires that $2^{\mathcal{O}}$ is totally ordered (e.g., by numbering the objects and using the lexicographical order for object sets). The algorithm enumerates object sets according to the lexicographical order, and applies C . The process starts with $C(\emptyset)$, which determines the extent of the bottom element. Once an extent has been found, its lexicographical successors are enumerated and C is applied, until the next extent (in lexicographic order) is found.

Construction of concept lattices and implication bases has typical time complexity $O(n^3)$ for an $n \times n$ table, but can be exponential in the worst case. Empirical studies show that even for large tables, exponential behavior is extremely rare [28]. In fact, it can be shown that if the number of attributes for every object is bounded (which is true for most applications), the lattice size is linear in the number of table entries [15]. In practice, Ganter’s algorithm needs less than a second for 2000-element lattices on a standard workstation [28].

If a row or column is added to a table, the lattice for the original table is a sublattice of the lattice for the extended table, and the new lattice can be constructed incrementally from the old one. The minimal implication base can be constructed in an incremental manner as well [13].

2.4 Algebraic decomposition

The structure theory of concept lattices allows for even deeper insight into the original relation. Without going into details, we just mention three important decomposition techniques:

- A *horizontal decomposition* is possible, if the lattice consists of independent sublattices connected only via the top and bottom element. Even if there are a few so-called *interferences* (infima between sublattices), a horizontal decomposition is still possible after interference removal.
- *Congruences* group lattice elements into classes such that the classes again form a lattice, thus grouping “related” elements into one congruence class. Concept lattices allow for a particular elegant characterization of congruences [13]. There is also a notion of “weak congruences”, called *block relations* where “congruences” are non-transitive and a congruence class corresponds to a rectangle shape in the table.

- *Subdirect decomposition* tries to construct the lattice as a (sublattice of) a cartesian product of two or more smaller lattices. Effective algorithms for subdirect decomposition exist [12].

There is much more to say about concept lattices, and related algorithms and methodology. Davey and Priestley’s book [10] contains a chapter on elementary concept analysis. Ganter and Wille [13] treat the topic in depth.

3 Objects and attributes

Roughly speaking, the objects and attributes in our domain are variables and class members, respectively, and the table that will be constructed in Section 4 identifies for each variable which members must be included in its type. Before we can define the objects and attributes more precisely, we need to introduce some terminology. In what follows, \mathcal{P} denotes a program containing a class hierarchy, or a collection of programs that share a class hierarchy. Further, v, w, \dots denote the variables in \mathcal{P} whose type is a class, and p, q, \dots the variables in \mathcal{P} whose type is a pointer to a class (references can be treated similarly, and we omit their formalization in the present paper). Expressions are denoted by x, y, \dots . We will henceforth use “variables” to refer to variables as well as parameters. In the definitions that follow, $TypeOf(\mathcal{P}, x)$ denotes the type of expression x in \mathcal{P} .

The *objects* of our domain are the program variables through which the class hierarchy is accessed. Variables whose type is (pointer to) built-in can be ignored because the class hierarchy can only be accessed through variables whose type is *class-related* (i.e., variables whose type is a class, or a pointer to a class). Definition 1 below defines sets of variables $ClassVars$ and $ClassPtrVars$ whose type is a class, and a pointer to a class, respectively. In Section 6.1, we will discuss how to model heap-allocated objects. Note that $ClassPtrVars$ includes implicitly declared `this` pointers of methods. In order to distinguish between `this` pointers of different methods, we will henceforth refer to the `this` pointer of method `A::f()` by the fully qualified name of its method, i.e., `A::f`.

Definition 1 *Let \mathcal{P} be a program. Then, the set of class-typed variables and the set of pointer-to-class-typed variables are defined as follows:*

$$\begin{aligned} ClassVars(\mathcal{P}) &\triangleq \\ &\{ v \mid v \text{ is a variable in } \mathcal{P}, TypeOf(\mathcal{P}, v) = C, \text{ for some class } C \text{ in } \mathcal{P} \} \\ ClassPtrVars(\mathcal{P}) &\triangleq \\ &\{ p \mid p \text{ is a variable in } \mathcal{P}, TypeOf(\mathcal{P}, *p) = C, \text{ for some class } C \text{ in } \mathcal{P} \} \end{aligned}$$

The *attributes* of our domain are class members. Following the definitions of [35, 36], we will distinguish between *definitions* and *declarations* of members. We define these terms as follows: The definition of a member comprises a member’s signature (interface) as well as the executable code in its body, whereas the declaration of a member only represents its signature. This distinction is needed for accurately modeling virtual method calls.

```

class A {
public:
    virtual int f(){ return g(); };
    virtual int g(){ return x; };
    int x;
};
class B : public A {
public:
    virtual int g(){ return y; };
    int y;
};
class C : public B {
public:
    virtual int f(){ return g() + z; };
    int z;
};

int main(){
    A a; B b; C c;
    A *ap;
    if (...) { ap = &a; }
    else { if (...) { ap = &b; }
           else { ap = &c; } }
    ap->f();
    return 0;
}

```

Figure 4: Example program \mathcal{P}_1 .

Consider a call to a virtual method f from a *pointer* p . In this case, only the declaration of f needs to be contained in p 's type in order to be able to invoke f ; the body of f does not need to be statically visible to p ⁵. Naturally, a *definition* of f must be visible to the object that p points to at run-time, so that the dynamic dispatch can be executed correctly.

Definition 2 (shown below) defines sets $MemberDcls(\mathcal{P})$ and $MemberDefs(\mathcal{P})$ of member declarations and member definitions in \mathcal{P} . We distinguish between declarations and definitions of virtual methods for the reasons stated above. For nonvirtual methods, making this distinction is not necessary because the full definition of a nonvirtual method must always be statically visible to the caller. Therefore, nonvirtual methods are modeled using definitions only. Data members are modeled as declarations because they have no `this` pointer from which other members can be accessed.

Definition 2 *Let \mathcal{P} be a program. Then, we define the set of member declarations and member definitions as follows:*

$$\begin{aligned}
 MemberDcls(\mathcal{P}) &\triangleq \\
 &\{ dcl(C::m) \mid m \text{ is a data member or virtual method in class } C \} \\
 MemberDefs(\mathcal{P}) &\triangleq \\
 &\{ def(C::m) \mid m \text{ is a virtual or nonvirtual method in class } C \}
 \end{aligned}$$

Example: Figure 4 shows a program \mathcal{P}_1 that will be used as a running example. For \mathcal{P}_1 , we have:

$$\begin{aligned}
 ClassVars(\mathcal{P}_1) &\equiv \{ a, b, c \} \\
 ClassPtrVars(\mathcal{P}_1) &\equiv \{ ap, A::f, A::g, B::g, C::f \} \\
 MemberDcls(\mathcal{P}_1) &\equiv \{ dcl(A::f), dcl(A::g), dcl(A::x), dcl(B::g), \\
 &\quad dcl(B::y), dcl(C::f), dcl(C::z) \} \\
 MemberDefs(\mathcal{P}_1) &\equiv \{ def(A::f), def(A::g), def(B::g), def(C::f) \}
 \end{aligned}$$

⁵Our objective is to identify the smallest possible set of member declarations and definitions that must be included in the type of any variable. Including the *definition* of f in $*p$'s type may lead to the incorporation of members that are otherwise not needed (in particular, members accessed from f 's `this` pointer).

In Section 6.2, we will discuss how class-typed data members (which behave like variables because other members can be accessed from them) are modeled.

4 Table Construction

This section describes how tables and lattices are constructed. Recall that the purpose of the table is to record for each variable the set of members that are used. A few auxiliary definitions will be presented first, in Section 4.1.

4.1 Auxiliary definitions

For each variable v in $ClassPtrVars(\mathcal{P})$ we will need a conservative approximation of the variables in $ClassVars(\mathcal{P})$ variables that v may point to. Any of several existing algorithms [8, 22, 31, 26] can be used to compute this information, and we do not make assumptions about the particular algorithm used to compute points-to information. Definition 3 expresses the information supplied by some points-to analysis algorithm as a set $PointsTo(\mathcal{P})$, which contains a pair $\langle p, v \rangle$ for each pointer p that may point to a class-typed variable v .

Definition 3 *Let \mathcal{P} be a program. Then, the points-to information for \mathcal{P} is defined as follows:*

$$PointsTo(\mathcal{P}) \triangleq \{ \langle p, v \rangle \mid p \in ClassPtrVars(\mathcal{P}), v \in ClassVars(\mathcal{P}), p \text{ may point to } v \}$$

Example: We will use the following points-to information for program \mathcal{P}_1 . Recall that $X::f$ denotes the `this` pointer of method $X::f()$.

$$PointsTo(\mathcal{P}_1) \equiv \{ \langle ap, a \rangle, \langle ap, b \rangle, \langle ap, c \rangle, \langle A::f, a \rangle, \langle A::f, b \rangle, \langle C::f, c \rangle, \langle A::g, a \rangle, \langle B::g, b \rangle, \langle B::g, c \rangle \}$$

Note that the following simple algorithm suffices to compute the information of Example 4.1: for each pointer p of type $*X$, assume that it may point to any object of type Y , such that (i) $Y = X$ or Y is a class transitively derived from X , and (ii) if p is the `this` pointer of a virtual method $C::m$, no overriding definitions of m are visible in class Y .

We will use the following terminology for function and method calls. A *direct* call is any call to a function or a nonvirtual method, or an invocation of a virtual method from a variable in $ClassVars(\mathcal{P})$. An *indirect* call is an invocation of a virtual method from a variable in $ClassPtrVars(\mathcal{P})$ (requiring a dynamic dispatch).

4.2 Table entries for member access operations

Table T has a *row* for each element of $ClassVars(\mathcal{P})$ and $ClassPtrVars(\mathcal{P})$, and a *column* for each element of $MemberDcls(\mathcal{P})$ and $MemberDefs(\mathcal{P})$. Informally, an entry $(y, dcl(A::m))$

appears in T iff the declaration of m is contained in y 's type, and an entry $(y, \text{def}(A::m))$ appears in T iff the definition of m is contained in y 's type. We begin by adding entries to T that reflect the member access operations in the program. Definition 4 below defines a set $\text{MemberAccess}(\mathcal{P})$ of all pairs $\langle m, y \rangle$ such that member m is accessed from variable y . For an *indirect* call $p \rightarrow f(y_1, \dots, y_n)$, we also include an element $\langle f, y \rangle$ in $\text{MemberAccess}(\mathcal{P})$ for each $\langle p, y \rangle \in \text{PointsTo}(\mathcal{P})$.

Definition 4 *Let \mathcal{P} be a program. Then, the set of member access operations in \mathcal{P} is defined as follows:*

$$\begin{aligned} \text{MemberAccess}(\mathcal{P}) \triangleq & \\ & \{ \langle m, v \rangle \mid v.m \text{ occurs in } \mathcal{P}, m \text{ is a class member in } \mathcal{P}, v \in \text{ClassVars}(\mathcal{P}) \} \cup \\ & \{ \langle m, *p \rangle \mid p \rightarrow m \text{ occurs in } \mathcal{P}, m \text{ is a class member in } \mathcal{P}, \\ & \quad p \in \text{ClassPtrVars}(\mathcal{P}) \} \cup \\ & \{ \langle m, y \rangle \mid p \rightarrow m \text{ occurs in } \mathcal{P}, \langle p, y \rangle \in \text{PointsTo}(\mathcal{P}), m \text{ is a virtual method in } \mathcal{P} \} \end{aligned}$$

Example: For program \mathcal{P}_1 of Figure 4, we have:

$$\begin{aligned} \text{MemberAccess}(\mathcal{P}_1) \equiv & \\ & \{ \langle x, *A::g \rangle, \langle y, *B::g \rangle, \langle z, *C::f \rangle, \langle g, *A::f \rangle, \langle g, *C::f \rangle, \\ & \quad \langle f, *ap \rangle, \langle f, a \rangle, \langle f, b \rangle, \langle f, c \rangle, \langle g, a \rangle, \langle g, b \rangle, \langle g, c \rangle \} \end{aligned}$$

Accessing a class member is not an entirely trivial operation because different classes in a class hierarchy may contain members with the same name (or signature). Furthermore, in the presence of multiple inheritance, an object may contain multiple subobjects of a given type C , and hence multiple members $C::m$. This implies that whenever a member m is accessed, one needs to determine which m is being selected. This selection process is defined informally in the C++ Draft Standard [1] as a set of rules that determine when a member *hides* or *dominates* another member with the same name. Rossie and Friedman [25] provided a formalization of the member lookup, as a function on *subobject graphs*. This framework has subsequently been used by Tip et al. as a formal basis for operations on class hierarchies such as slicing [33] and specialization [35]. Ramalingam and Srinivasan recently presented an efficient algorithm for member lookup [24].

For the purposes of the present paper, we will assume the availability of a function *static-lookup* which, given a class C and a member m , determines the base class B (B is either C , or a transitive base class of C) in which the selected member is located⁶. For details on function *static-lookup*, the reader is referred to [25, 33].

We are now in a position to state how the appropriate relations between variables and declarations and definitions should be added to the table:

⁶In [25, 33], *static-lookup* is defined as a function from subobject to subobjects. Since the present paper is only concerned with the *classes* in which members are located, we will simply ignore all subobject information below.

Definition 5 Let \mathcal{P} be a program with associated table T . Then, the following entries are added to the table due to member access operations that occur in the program.

$$\frac{\langle m, y \rangle \in \text{MemberAccess}(\mathcal{P}), \quad m \in \text{DataMembers}(\mathcal{P}), \\ X \equiv \text{static-lookup}(\text{TypeOf}(\mathcal{P}, y), m)}{(y, \text{dcl}(X::m)) \in T}$$

$$\frac{\langle m, y \rangle \in \text{MemberAccess}(\mathcal{P}), \quad m \in \text{NonVirtualMethods}(\mathcal{P}), \\ X \equiv \text{static-lookup}(\text{TypeOf}(\mathcal{P}, y), m)}{(y, \text{def}(X::m)) \in T}$$

$$\frac{\langle m, y \rangle \in \text{MemberAccess}(\mathcal{P}), \quad m \in \text{VirtualMethods}(\mathcal{P}), \\ y \equiv *p, \quad p \in \text{ClassPtrVars}(\mathcal{P}), \\ X \equiv \text{static-lookup}(\text{TypeOf}(\mathcal{P}, y), m)}{(y, \text{dcl}(X::m)) \in T}$$

$$\frac{\langle m, y \rangle \in \text{MemberAccess}(\mathcal{P}), \quad m \in \text{VirtualMethods}(\mathcal{P}), \\ y \equiv v, \quad v \in \text{ClassVars}(\mathcal{P}), \\ X \equiv \text{static-lookup}(\text{TypeOf}(\mathcal{P}, y), m)}{(y, \text{def}(X::m)) \in T}$$

4.3 Table entries for this pointers

The next table construction rule we will present is concerned with **this** pointers of methods. Consider the fact that for each method $C::f()$, there is a column in the table labeled $\text{def}(C::f)$, and a row labeled $*C::f$. The former is used to express the fact that method $C::f()$ may be called from objects. The latter is necessary to reflect members being accessed from method $C::f()$'s **this** pointer. Unless precautions are taken, the attribute $\text{def}(C::f)$ and the object $*C::f$ may appear at different points in the lattice, though $\gamma(*C::f) \geq \mu(\text{def}(C::f))$ must always hold⁷. In such cases, our method effectively infers that the type of a **this** pointer could be a base class of the type in which method $C::f$ occurs (and therefore be less constrained). However, in reality, the type of a method's **this** pointer is *determined by* the class in which the associated method definition appears.

The table entries added by Definition 6 will force a method's attribute and a method's **this** pointer to appear at the same lattice element; by ensuring $\gamma(*C::f) \leq \mu(\text{def}(C::f))$. This will allow us later to remove rows for **this** pointers from the table when constructing the lattice.

Definition 6 Let \mathcal{P} be a program. Then, the following entries are added to the table:

$$\frac{\text{def}(C::m) \in \text{MemberDefs}(\mathcal{P})}{(*C::m, \text{def}(C::m)) \in T}$$

⁷See Appendix.

	dc1(A::f)	dc1(A::g)	dc1(A::x)	def(A::f)	def(A::g)	dc1(B::g)	dc1(B::Y)	def(B::g)	dc1(C::z)	def(C::f)
a				×	×					
b				×				×		
c								×		×
*ap	×									
*A::f		×		×						
*A::g			×		×					
*B::g							×	×		
*C::f						×			×	×

Table 1: Initial table for program \mathcal{P}_1 of Figure 4. Arrows indicate implications due to assignments (see Section 4.4).

Example: Table 1 shows the table for program \mathcal{P}_1 of Figure 4 after adding the entries according to Definitions 5 and 6. The purpose of the arrows at the side of the table will be explained in Section 4.4.

4.4 Table entries for assignments

Consider an assignment $x = y$, where $x \equiv v$ and $y \equiv w$, for some class-typed variables $v, w \in \text{ClassVars}(\mathcal{P})$. Such an assignment is only valid if the type of x is a base class of the type of y . Consequently, any member declaration or definition that occurs in x 's type must also occur in y 's type. We will enforce this constraint using an *implication* from the row for x to the row for y . However, we will begin by formalizing the notion of an assignment.

Definition 7 below defines a set $\text{Assignments}(\mathcal{P})$ that contains a pair of objects $\langle v, w \rangle$ for each assignment $v = w$ in \mathcal{P} where v and w are class-typed. In addition, $\text{Assignments}(\mathcal{P})$ also contains entries for cases where the type of the left-hand side and/or the right-hand side of the assignment are a pointer to a class. Parameter-passing in direct calls to functions and methods is modeled by way of assignments between corresponding formal and actual parameters. For an *indirect* call $p \rightarrow f(y_1, \dots, y_n)$, $\text{Assignments}(\mathcal{P})$ contains additional elements that model the parameter-passing in the direct call $x.f(y_1, \dots, y_n)$, for each $\langle p, x \rangle \in \text{PointsTo}(\mathcal{P})$. That is, we conservatively approximate the potential targets of dynamically dispatched calls. The set $\text{Assignments}(\mathcal{P})$ will also contain elements for implicit parameters such as **this** pointers of methods and function/method return values whose type is class-related.

Definition 7 *Let \mathcal{P} be a program. Then, the set of assignments between variables whose*

type is a (pointer to a) class is defined as follows:

$$\begin{aligned} \text{Assignments}(\mathcal{P}) \triangleq & \\ & \{ \langle v, w \rangle \mid v = w \text{ occurs in } \mathcal{P}, v, w \in \text{ClassVars}(\mathcal{P}) \} \cup \\ & \{ \langle *p, w \rangle \mid p = \&w \text{ occurs in } \mathcal{P}, p \in \text{ClassPtrVars}(\mathcal{P}), w \in \text{ClassVars}(\mathcal{P}) \} \cup \\ & \{ \langle *p, *q \rangle \mid p = q \text{ occurs in } \mathcal{P}, p, q \in \text{ClassPtrVars}(\mathcal{P}) \} \cup \\ & \{ \langle *p, w \rangle \mid *p = w \text{ occurs in } \mathcal{P}, p \in \text{ClassPtrVars}(\mathcal{P}), w \in \text{ClassVars}(\mathcal{P}) \} \cup \\ & \{ \langle v, *q \rangle \mid v = *q \text{ occurs in } \mathcal{P}, v \in \text{ClassVars}(\mathcal{P}), q \in \text{ClassPtrVars}(\mathcal{P}) \} \cup \\ & \{ \langle *p, *q \rangle \mid *p = *q \text{ occurs in } \mathcal{P}, p, q \in \text{ClassPtrVars}(\mathcal{P}) \} \end{aligned}$$

Example: For program \mathcal{P}_1 of Figure 4, we have:

$$\begin{aligned} \text{Assignments}(\mathcal{P}_1) \equiv & \\ & \{ \langle *ap, a \rangle, \langle *ap, b \rangle, \langle *ap, c \rangle, \langle *A::f, a \rangle, \langle *A::f, b \rangle, \\ & \quad \langle *C::f, c \rangle, \langle *A::g, a \rangle, \langle *B::g, b \rangle, \langle *B::g, c \rangle \} \end{aligned}$$

We are now in a position to express how elements should be added to the table due to assignments. Definition 8 states this as an *implication*, which tells us how elements should be copied from one row to another.

Definition 8 *Let \mathcal{P} be a program with associated table T . Then, the following implications must be encoded in the table due to assignments that occur in \mathcal{P} :*

$$\frac{\langle x, y \rangle \in \text{Assignments}(\mathcal{P})}{x \rightarrow y}$$

Note that assignment implications are implications between “objects” (in the sense of concept analysis), hence an assignment implication $x \rightarrow y$ causes x to appear above y (i.e. $\gamma(x) \geq \gamma(y)$) in the lattice. Cyclic assignments will generate cyclic implications, which will collapse the corresponding lattice elements into one point: all the involved variables must have the same type.

Example: For program \mathcal{P}_1 of Figure 4, the following assignment implications are generated:

$$\begin{aligned} *ap \rightarrow a, *ap \rightarrow b, *ap \rightarrow c, *A::f \rightarrow a, *A::f \rightarrow b, \\ *C::f \rightarrow c, *A::g \rightarrow a, *B::g \rightarrow b, *B::g \rightarrow c \end{aligned}$$

These implications are indicated on the left side of Table 1. Table 2 is obtained by copying the elements from the “source row” to the “target row” according to each of these implications.

4.5 Table entries for preserving dominance/hiding

The table thus far encodes for each variable the members contained in its type (either directly because a member is accessed from that variable, or indirectly due to assignments between variables). However, in the original class hierarchy, an object’s type may contain

	dcl(A::f)	dcl(A::g)	dcl(A::x)	def(A::f)	def(A::g)	dcl(B::g)	dcl(B::y)	def(B::g)	dcl(C::z)	def(C::f)
a	×	×	×	×	×					
b	×	×		×			×	×		
c	×					×	×	×	×	×
*ap	×									
*A::f		×		×						
*A::g			×		×					
*B::g							×	×		
*C::f						×			×	×

Table 2: Table after application of assignment implications. Arrows indicate implications for preserving hiding/dominance among members with the same name (see Section 4.5).

more than one member with a given name. In such cases, the member lookup rules of [1] determine which member is accessed. This is expressed as a set of rules that determine when a member *hides* or *dominates* another member with the same name. In cases where a variable contains two members m that have a hiding relationship in the original class hierarchy, this hiding relationship must be preserved: we are interested in generating a restructured hierarchy from the table, and the member access operations in the program might otherwise become ambiguous. Definition 9 incorporates the appropriate hiding/dominance relations into the table, using implications between attributes:

Definition 9 Let \mathcal{P} be a program with associated table T . Then, the following implications are incorporated into T in order to preserve hiding and dominance:

$$\frac{(x, \text{dcl}(A::m)) \in T, (x, \text{dcl}(B::m)) \in T, \text{ } A \text{ is a transitive base class of } B}{\text{dcl}(B::m) \rightarrow \text{dcl}(A::m)}$$

$$\frac{(x, \text{dcl}(A::m)) \in T, (x, \text{def}(B::m)) \in T, \text{ } A = B \text{ or } A \text{ is a transitive base class of } B}{\text{def}(B::m) \rightarrow \text{dcl}(A::m)}$$

$$\frac{(x, \text{def}(A::m)) \in T, (x, \text{def}(B::m)) \in T, \text{ } A \text{ is a transitive base class of } B}{\text{def}(B::m) \rightarrow \text{def}(A::m)}$$

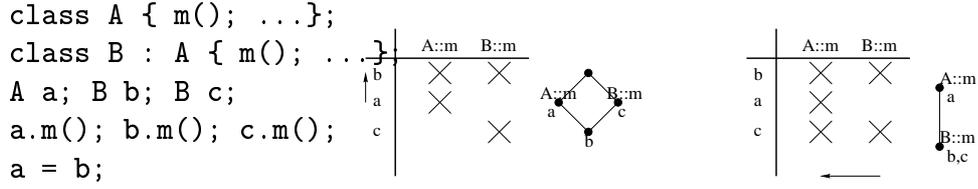


Figure 5: Effect of dominance rules

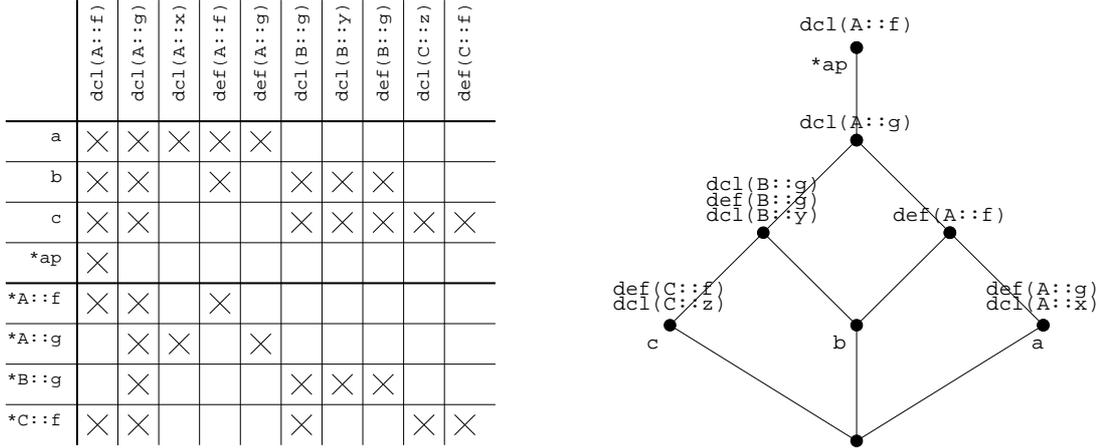


Figure 6: Final table and lattice for program \mathcal{P}_1 , after removing the rows labeled $*A::f$, $*A::g$, $*B::g$, and $*C::f$.

$$\frac{(x, \text{def}(A::m)) \in T, (x, \text{dcl}(B::m)) \in T, \quad A \text{ is a transitive base class of } B}{\text{dcl}(B::m) \rightarrow \text{def}(A::m)}$$

Dominance implications are implications between “attributes” (in the sense of concept analysis), hence a dominance implication $B::m \rightarrow A::m$ will cause $B::m$ to appear below $A::m$ (i.e. $\mu(B::m) \leq \mu(A::m)$) in the lattice. Due to the condition “ A is a (transitive) base class of B ”, dominance implications always connect subclass members to superclass members and cannot contain cycles (if $A = B$, a one-point “cycle” is generated).

Note the symmetry between dominance implications and assignment implications: the latter are implications between columns (attributes) and serve to preserve behavior of member lookup; the latter are implications between rows (objects) and serve to preserve behavior of subobject selection. Figure 5 demonstrates the effect of the dominance rules: subclass B of class A redefines method m . In the table, the implication due to assignment $a = b$; forces row a to be added to row b . But now the member access $b.m()$ has become ambiguous: the row for b contains entries for both $A::m$ and $B::m$. According to the dominance rules, an implication $B::m \rightarrow A::m$ is generated, which adds entry $(b, A::m)$ to the table. The corresponding lattice is a two-element chain and thus reproduces the original hierarchy, thereby reestablishing the dominance of $B::m$ over $A::m$.

Example: For program \mathcal{P}_1 , the following dominance implications are generated:

$$\begin{array}{llll} \text{def}(\mathbf{A}::\mathbf{f}) \rightarrow \text{dcl}(\mathbf{A}::\mathbf{f}) & \text{def}(\mathbf{A}::\mathbf{g}) \rightarrow \text{dcl}(\mathbf{A}::\mathbf{g}) & \text{def}(\mathbf{B}::\mathbf{g}) \rightarrow \text{dcl}(\mathbf{A}::\mathbf{g}) & \\ \text{dcl}(\mathbf{B}::\mathbf{g}) \rightarrow \text{dcl}(\mathbf{A}::\mathbf{g}) & \text{def}(\mathbf{B}::\mathbf{g}) \rightarrow \text{dcl}(\mathbf{B}::\mathbf{g}) & \text{def}(\mathbf{C}::\mathbf{f}) \rightarrow \text{dcl}(\mathbf{A}::\mathbf{f}) & \end{array}$$

These implications are shown at the bottom of Table 2. After incorporating these implications, the table in Figure 6 results.

Remark: Observe that the implication $\text{def}(\mathbf{B}::\mathbf{g}) \rightarrow \text{dcl}(\mathbf{A}::\mathbf{g})$ only becomes necessary after propagating table elements according to the other implications.

5 The new Hierarchy

5.1 Lattice construction

Since the assignment implications can generate new dominance implications and vice versa, a fixpoint iteration is necessary in order to compute the final table. This algorithm is described in section 8. After the table has converged, the lattice is constructed using Ganter’s algorithm (see section 2). As explained above, it can be directly interpreted as a new class hierarchy.

There is one issue concerning pointers that deserves mentioning. Recall that in Section 4.3 table entries were added to ensure that method definitions and their `this` pointers show up at the same lattice element. In order to avoid presenting redundant information to the user, we will henceforth omit `this` pointers from the lattice. The easiest way to accomplish this is to remove the rows for `this` pointer variables to the table prior to generating the lattice. Note that rows for `this` pointers cannot be left out during table construction because they are needed to model member accesses from `this` pointers, and the elements in such rows may be involved in implications due to assignments and dominance.

Example: Figure 6 shows the lattice for program \mathcal{P}_1 , generated from the final table after removing the rows labeled `*A::f`, `*A::g`, `*B::g`, and `*C::f`. The lattice can be interpreted directly as a new class hierarchy. It demonstrates that `a` does not access `B::y` and `C::z`, while `b` and `c` do not access `A::x` and `b` does not access `C::z`. Similarly, the lattice shows the fine-grained differences in method access: for example, `c` does not need `def(A::f)` and `def(A::g)`. Thus `a`, `b`, `c` will receive new types. The program statements are unchanged, but according to the new hierarchy, both `b` and `c` have become smaller.

Note that from a space optimization viewpoint, the lattice can be simplified further: for example, the two topmost elements could be merged (as they only contain method declarations), and even the edge $b \leftrightarrow \text{def}(\mathbf{A}::\mathbf{f})$ could be merged with the parallel edge. Tip and Sweeney [36] discuss such “peephole optimizations” in detail.

5.2 Properties of the lattice

The lattice, being a concept lattice, enjoys several important properties [13].

- The lattice is the smallest lattice compatible with the table and thus can be order-embedded into any other lattice compatible with the table. In fact, if the table represents a partial order, then the lattice is the Dedekind-McNeill completion of this partial order. Hence, the lattice is *minimal*.
- Attribute labels of lattice elements always occur as far upward in the lattice as possible (see section 2). Since attribute labels correspond to members in the classes of the new hierarchy, common members are factored out as much as possible. The same applies to common variables, which are factored out downwards as much as possible. Therefore the lattice is *maximally factorized*.

More important than minimality and maximal factorization are semantic properties of the lattice. In [36] it was proved that the assignment constraints and the dominance constraints guarantee

- *preservation of assignment behavior*: every assignment will select the same subobject from the right-hand-side object as in the original program;
- *preservation of lookup behavior*: every method call will select the same method definition via dynamic lookup as in the original program.

The lattice, interpreted as a new class hierarchy, respects all assignment and dominance constraints by construction. Since the statements of the program are unchanged, we thus can guarantee that *the new hierarchy is operationally equivalent to the old one*.

The lattice may contain elements which are neither labelled with an attribute nor an object (e.g. the center element in figure 3). Such elements are called “empty” and serve merely to group the members of other classes. In our application, an empty element C corresponds to a class which neither has any members, nor does any variable have type C in the new hierarchy. Section 7 will explain how to get rid of empty elements.

Remember that rows for **this**-pointer have been removed from the final table without semantic effect. A similar simplification can be applied to pointers in general. The lattice may be very fine-grained due to access patterns of pointers which basically have the same type, but access different members. Since a pointer will always appear above any object it may point to (see proof in the appendix), rows for pointers can safely be removed from the final table. The pointers can then be given the same type as the objects they point to – which still guarantees operational equivalence, since any pointer may still access all members it needs⁸. Of course, pointer rows are essential during table generation, as explained above for the special case of **this** pointers.

How does the final lattice depend on the precision of the points-to analysis? Since any points-to analysis computes a conservative approximation, different analyses differ only with respect to the row entries they generate for pointer variables: the more precise the points-to analysis, the less entries any pointer row will have. That is, if table T_1 has been

⁸More precisely, the pointer is given the supremum type of all object types it may point to; this supremum also exists in the reduced lattice and is still below the pointer’s type in the full lattice.

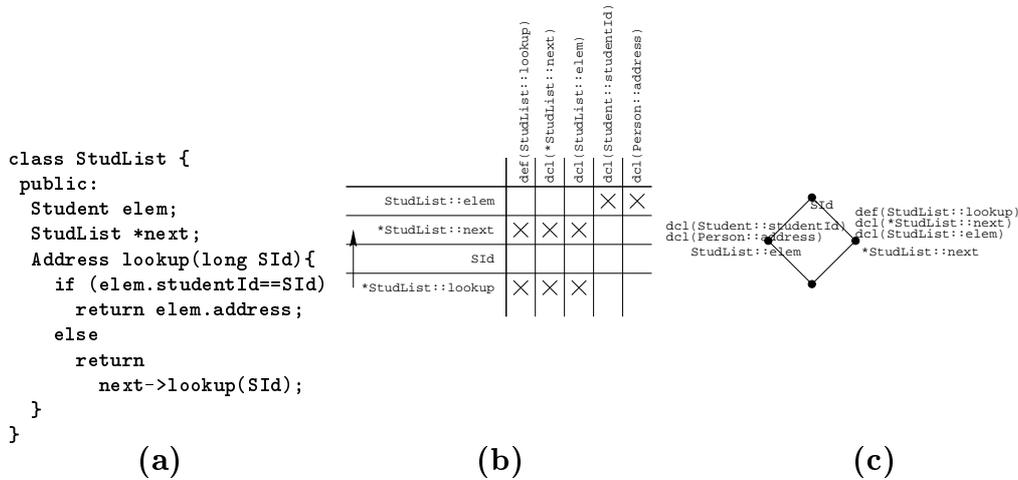


Figure 7: Analyzing a linked list of students

generated using a more precise points-to analysis than for table T_2 , we have $(o, a) \in T_1 \Rightarrow (o, a) \in T_2$. By the fundamental property $(o, a) \in T_{1,2} \iff \gamma_{1,2}(o) \leq \mu_{1,2}(a)$, hence $\gamma_1(o) \leq \mu_1(a) \Rightarrow \gamma_2(o) \leq \mu_2(a)$: $\mathcal{L}(T_1)$ can be order-embedded into $\mathcal{L}(T_2)$. In $T_{1,2}$, any possible pointer row is a superset of a “minimal row” for that pointer, which corresponds to the (undecidable) limit case of precise points-to analysis. By the above embedding, the lattice for the limit case can be found as a substructure in every actual lattice.

6 Language Details

The basic process for constructing tables and lattices, as described in the previous section, did not address a number of language features, that are not core issues, but indispensable in practice. This section addresses a number of such issues.

6.1 Heap allocation

We handle heap-allocated objects in a straightforward way by simply treating each allocation site in the program as a class-typed variable (e.g., an element of the set *ClassVars*). For the program of Figure 1, there are four such allocation sites, which we refer to as **Student1**, **Student2**, **Professor1**, and **Professor2**. In principle, more sophisticated context-sensitive analyses could be used to distinguish heap-allocated objects in different calling contexts, but we expect the benefits of this additional precision to be limited.

6.2 Modeling nested objects

The treatment of *class-related data members* (i.e., data members whose type is class-related such as **Student::advisor** in Figure 1) is an important issue. Like data members of built-in types, class-related data members can be accessed from variables and are therefore

modeled as attributes. However, since other members may be accessed from a class-related data member, such data members play an alternate role as objects.

In order to clarify the issues involved in the reengineering of such “nested” structures, consider a class C that contains a data member m whose type is some class D . Then, the following information about m is made explicit in the concept lattice:

- The set of variables in which m is contained. This is modeled by treating m as an “attribute” a . Any object that occurs below a in the lattice contains m .
- The set of members contained in the type of m . This is modeled by treating the type of m as an “object” o . The set of members contained in o correspond to the attributes that occur above o in the lattice. This set of members is a subset of the members of D in the original class hierarchy.

Note that the “attribute view” of m corresponds exactly to the way we previously modeled data members with a built-in type, whereas the “object view” of m corresponds exactly to the way we previously modeled variables. The definitions that are concerned with variables therefore apply to class-related data members as well, and for convenience we will henceforth assume the term “variable” to include class-related data members.

Figure 7(a) shows an example program that illustrates the issues related to class-related data members. Here, the example of Figure 1 is extended with a linked list of students. Observe that the data members `address` and `studentId` are accessed from the `elem` member of a student list, and that the `lookup` method of class `StudList` is accessed from the `next` data member. The data members `elem` and `next` are accessed from method `lookup`’s `this` pointer, and note that as usual, a table entry is added for the associated method definition of `lookup`. Furthermore, there is an assignment implication `*StudList::lookup = next` due to the recursive call.

Figure 7(b) shows the table after applying this implication, and Figure 7(c) shows the associated lattice. Let us call the left element $Student'$ and the right element $StudList'$, as suggested by their labels. Both `next` resp. `elem` show up twice in the lattice. $Student'$ shows that the *type* of `StudList::elem` must be $Student'$; $StudList'$ shows that the *member* `StudList::elem` must be located within class $StudList'$. Similar observations are valid for `next`: it must be located in $StudList'$ and also has type $StudList'$ – as to be expected. The example thus illustrates the dual role of class-typed data members: the lattice will not only display their type, but also their position in the class hierarchy.

6.3 Modeling constructors

Constructors require special attention. A constructor generally initializes all data members contained in an object. If no constructor is provided by the user, a so-called default constructor is generated by the compiler, which performs the necessary initializations. The compiler may also generate a *call* to a constructor in certain cases. Modeling these compiler-generated actions as member access operations would lead us to believe that each member m of class C is needed in all C -instances, even in cases where the only access to m

<pre> class O { ... }; class A extends O { ... }; class B extends O { ... }; class Example { public static void main(String args[]){ A a = new A(); O o = a; if (o instanceof A){ /* reached */ } if (o instanceof B){ /* unreached */ } A a2 = (A)o; /* succeeds */ A a3 = (B)o; /* ClassCastException */ } } </pre>	<pre> class O { boolean isA(){ return false; }; boolean isB(){ return false; }; A toA(){ throw new ClassCastException(); }; B toB(){ throw new ClassCastException(); }; }; class A extends O { boolean isA(){ return true; }; A toA(){ return this; }; ... }; class B extends O { boolean isB(){ return true; }; B toB(){ return this; }; ... }; class Example { public static void main(String args[]){ A a = new A(); O o = a; if (o.isA()){ /* reached */ } if (o.isB()){ /* unreached */ } A a2 = o.toA(); /* succeeds */ B b = o.toB(); /* ClassCastException */ } } </pre>
(a)	(b)

Figure 8: (a) Example Java program that uses type cast and `instanceof` operations. (b) Equivalent Java program after transforming away `instanceof` and cast operations.

consists of its (default) initialization. Compiler-generated constructors, compiler-generated initializations, and compiler-generated calls to constructors will therefore be excluded from the set of member access operations. Destructors can be handled similarly.

6.4 Type Casts and Instanceof operations

Modern object-oriented languages such as Java provide language features for testing the run-time type of an object, or down-casting an object to a derived type. Since these operations are used heavily, any realistic implementation will need to deal with them.

Our approach to dealing with type cast and instance-of operations will be to transform them into a semantically equivalent piece of code consisting of only virtual method calls and exception-handling constructs. We will outline these transformations for the cast and `instanceof` operations as they are used in Java. In Java, these operations have the following semantics:

- An expression `e instanceof C` evaluates to `true` if the run-time type of the object pointed to by reference `e` is `C` or a subclass of `C`. Otherwise, the expression evaluates to `false`.
- A cast-expression `(C)e` evaluates to an expression `e` with static type `C` if the run-time type of `e` is `C` or a subclass of `C`. Otherwise, an exception of type `ClassCastException`

is thrown.

Our strategy for transforming `instanceof`-expressions will be as follows⁹ For each type `C`, we introduce a method `isC()` in the root class `O` of the hierarchy. This method has return type `boolean`, and the default definition of `isC` in class `O` returns `false`. Class `C` provides an overriding definition¹⁰ of `isC()` that returns `true`. Now, every expression of the form `e instanceof C` is transformed into `e.isC()`. One can see easily that the expressions `e instanceof C` and `e.isC()` return `true` under exactly the same conditions.

Cast expressions are transformed in a similar manner. For each type `C`, we introduce a method `toC()` in the root class `O` of the hierarchy. This method has return type `C`, and the default definition of `isC` in class `O` throws an exception of type `ClassCastException`. Class `C` provides an overriding definition of `toC()` that returns `this`. Now, every expression of the form `(C)e` is transformed into `e.toC()`. It can easily be seen that the expressions `(C)e` and `e.toC()` succeed under exactly the same conditions.

Figure 8(a) shows an example program containing various `instanceof` and down-cast expressions. Figure 8(b) shows the program after transforming away all these expressions. After eliminating all cast and `instanceof` expressions, the resulting program can be processed with the techniques presented in the previous section.

After generating the lattice, the artificial `isC()` and `toC()` methods can easily be transformed back into cast and `instanceof` operations if the reengineer desires to do so. As an example, we will outline how `toC()` methods can be transformed back into cast operations in the transformed class hierarchy. Let x be the lattice element labeled `def(C.toC())`, and suppose that class name X has been associated with this lattice element. Then `e.toC()` can be transformed into `(X)e`.

6.5 Exceptions

Exception-handling constructs give rise to additional control flow, but do not influence member access patterns. Therefore, exceptions do not require special treatment. Note however that the abundance of (implicit or explicit) exceptions in Java complicates points-to analysis, and one might think of adopting factored control flow graphs [9] for our analysis.

6.6 Arrays

Arrays are treated as monolithic variables, and we do not distinguish between different array elements. One might think of integrating a fine-grained array analysis, such as the Omega test [23].

However, fine-grained array analysis may also reduce analysis precision in case of mixed co- and contravariance. In Java, arrays are covariant: $A \leq B$ implies $A[] \leq B[]$. If Java would allow fully contravariant method overriding (which it does not), we would have

⁹This transformation was proposed by M. Streckenbach.

¹⁰In cases where the target type C is an interface, this overriding definition should not be placed in C itself, but in all classes that implement C .

$dom(f_A) \geq dom(f_B)$ for any B -method f redefined in A . KABA translates array accesses into predefined method calls $a.store(i, x)$ and $a.access(i)$.¹¹ Thus, by contravariance for $store$'s second argument we would obtain $A = dom(store_{A[]})[2] \geq dom(store_{B[]})[2] = B$. Hence we suddenly have $A \leq B$ as well as $B \leq A$, collapsing lattice elements and hiding fine-grained access patterns.

6.7 Dynamic Class Loading

Java offers a mechanism for dynamic loading of classes, and also some reflective devices. For example, one might compute the name of a class at runtime, and then create an object of this class. Naturally, our analysis cannot handle such features; it has to make worst-case assumptions. In order to avoid massive loss of analysis precision another option is to rely on additional user input as in [34]: the reengineer might supply background knowledge about the expected outcome of dynamic constructs. Of course, if the reengineer's assumptions are false, operational equivalence is lost.

6.8 Multiple subobjects

If an object x contains multiple subobjects of some type C (due to the use of nonvirtual multiple inheritance), our tables do not make a distinction between the various "copies" of the members of C in x . This leads to problems if the objective is to generate a new hierarchy from the lattice in which the distinct copies of the members of C must be preserved. We consider this to be a minor problem because situations where nonvirtual inheritance is used for its "member replicating" effect are quite rare in practice, and the restructuring tool could inform the user of the cases where the problem occurs. A clean solution to this problem would involve the encoding of subobject information in the table using an adaptation of the approach of [35, 36].

7 Restructuring class hierarchies

7.1 Students and Professors reconsidered

Table 3 shows the final table for the example of Figure 1, as obtained by analyzing the class hierarchy along with the two example programs. The lattice corresponding to this table was shown previously in Figure 2 (note that we replaced member definitions by the corresponding method names there for convenience).

The following can be learned from the lattice:

- Data members that are not accessed anywhere in the program (e.g., `Person::socialSecurityNumber`) appear at the bottom element of the lattice.

¹¹This transformation is similar to the transformations for typecasts and instanceof expressions described above.

	decl(Person::name)	decl(Person:address)	decl(Person::socialSecurityNumber)	decl(Student::studentId)	decl(Student::advisor)	decl(Professor::faculty)	decl(Professor::workAddress)	decl(Professor::assistant)	def(Student::Student)	def(Student::setAdvisor)	def(Professor::Professor)	def(Professor::hireAssistant)
*s1				×					×			
*s2												
*p1												
*p2							×					×
*s												
*p												
Student1	×	×		×	×				×	×		
Student2	×	×		×					×			
Professor1	×					×	×	×			×	
Professor2	×					×	×	×			×	×
*advisor												
*assistant												
*Student::Student	×	×		×					×			
*Student::setAdvisor					×					×		
*Professor::Professor	×					×	×	×			×	
*Professor::hireAssistant								×				×

Table 3: Final table for the Student/Professor example.

- Data members of a base class B that are not used by (instances of) all derived classes of B are revealed. Such data members (e.g., `Person::address`) appear above (variables of) some but not all derived classes of B . For example, `Person::address` appears above instances of `Student`, but not above any instances of `Professor`.
- Variables from which no members are accessed appear at the top element of the lattice (e.g., `s`).
- Data members that are properly initialized appear above the (constructor) method that is supposed to initialize them. If this is not the case, the data member may not be initialized. For example, we know that `Student::Student` does not initialize `Student::advisor` because that data member does not appear above `Student::Student` in the lattice.
- Situations where instances of a given type C access different subsets of C 's members are revealed by the fact that variables of type C appear at different points in the lattice. Our example contains two examples of this phenomenon. The instances `Professor1` and `Professor2` of type `Professor` and the instances `Student1` and `Student2` of type `Student`.

As we mentioned earlier, a class hierarchy may be analyzed along with any number of programs, or without any program at all. The latter case may provide insights into the “internal structure” of a class library. Figure 9 shows the lattice obtained by analyzing the class hierarchy of Figure 1(a) *without* the programs of Figure 1(b) and (c); only code in method bodies is analyzed. Clearly, the resulting lattice should not be interpreted as a restructuring proposal, because it does not reflect the *usage* of the class hierarchy. However, there are some interesting things to note. For example, `socialSecurityNumber` is not accessed anywhere. If we would know in addition that `socialSecurityNumber` is private (i.e., that it can only be accessed by methods within its class), we could inform the user that it is effectively dead. Observe also that no members are accessed from method parameters `s` and `p`. Since the scope of these variables is local to the library, we know that analyzing additional code will not change this situation.

7.2 Restructuring transformations

Once the lattice has been computed and displayed, it can help in understanding the actual behaviour of a class hierarchy and thus serve as a basis for restructuring tasks (see section 9). In addition, several global restructuring transformations are possible:

- Unlabeled (“empty”) lattice elements correspond to classes without members and without variables using them. The lattice can be simplified by pruning all such elements, and directly connecting their subordinate and superordinate neighbors. The resulting structure is not a lattice anymore, but only a partial order, but this is not so important: a class hierarchy need only be a partial order, and lookup behavior

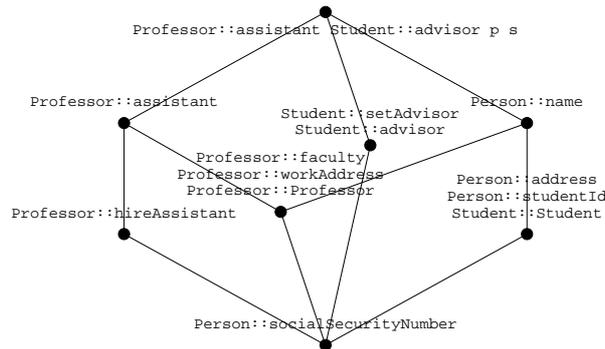


Figure 9: Lattice obtained by analyzing the class hierarchy of Figure 1 without accompanying programs.

and subobjects are not affected. Our case studies show that between 5 and 20 percent of all lattice elements are empty.

- A reduced lattice can be shown which contains only real objects, but no pointers. This lattice is obtained by deleting *all* pointer rows (not just **this**-pointers) from the final table. For reengineering purposes, this lattice seems more appropriate than the fine-grained one: ultimately, objects access members and determine the optimal class structure; fine-grained behavior of pointers is not helpful in an overall picture.¹² The resulting lattice is, by the theory of concept lattices, a sublattice of the original one, and still operationally equivalent. Our case studies show that up to 50 percent of the lattice elements disappear after removing pointer rows.
- The user can decide to merge adjacent lattice elements if the distinction between these concepts is irrelevant (possibly because the lattice reflects a specific use of the hierarchy). For example, one may decide that the distinction between professors that hire assistants, and professors that do not hire assistants is irrelevant, and therefore merge the concepts for **Professor1** and **Professor2**. There are however some issues that a tool must take into account, because we want to preserve member lookup behavior. For example, merging two concepts that have different definitions of a virtual method f associated with them is not possible, because at most one f can occur in any class. In general, merging must respect the dominance constraints for members.
- With certain limitations, the user may move attributes upwards in the lattice, and object downwards. For example, the user may decide that **socialSecurityNumber** should be retained in the restructured class hierarchy, and move the corresponding attribute up to the concept labeled with attribute **Person::name**. Again, dominance constraints must be respected.

¹²Remember that by construction, pointers always appear above the objects they point to.

- Background knowledge that is not reflected in the lattice, e.g., “the type of x must be a base class of the type of y ”, can be integrated via background implications. Technically, background implications are treated the same way as dominance implications.
- Color should be used to display relevant substructures in the lattice, e.g., display all variables which formerly had the same type, or all members which formerly were in the same class.
- Associations in the sense of OMT can be recovered from the occurrences of class-typed members, and corresponding arcs added to the lattice.
- For very large class hierarchies, the tool could allow the user to focus on a selected subhierarchy either by specifying its minimal and maximal elements in the lattice, or by selecting specific rows and columns in the table (e.g. those belonging to a specific class).
- The structure theory of concept lattices offers several algebraic decompositions, such as horizontal decomposition, interference analysis or block relations (see section 2.4). They can be used to measure quality factors such as cohesion and coupling [29, 18].
- Eventually, the user may associate names with lattice elements. When the reengineer is done manipulating the lattice, these names could be used as class names in the restructured hierarchy. For example, by examining the lattice, the programmer may determine that `Student` objects on which the `setAdvisor` method is invoked are graduate students, whereas `Student` objects on which this method is not called are undergraduates. Consequently, he may decide to associate names `Student` and `GraduateStudent` with the concepts labeled `s2` and `s1`, respectively.
- Finally, source code can be generated according to the new hierarchy; thereby utilizing the reduced object memory requirements and the improved structure in the new hierarchy.

7.3 Dealing with multiple inheritance

The analysis results are presented in form of a lattice, hence will naturally contain multiple inheritance if interpreted as a class hierarchy. Since Java does not support multiple inheritance, generated hierarchies may not be representable in Java source code. Note that if the meet point and its superclasses are in fact interface classes, there is no representation problem. Note further that multiple inheritance is only a problem if the method is to be used as a program transformation, but not if the lattice only serves program understanding.

Introducing a certain loss of precision, multiple inheritance can be removed as follows. Every occurrence of multiple inheritance leads to a “diamond” structure in the lattice, such as the diamond $Professor2 - Professor1 - p2 - Professor :: assistant$ in Figure 2. By moving members up and variables down (as explained above) the diamond can be

transformed into a simple chain, while still maintaining behavioral equivalence. In Figure 2, *p2* can be moved down to *Professor2*, while *Professor :: hireAssistant* can be moved up to *Professor :: assistant*. Finally the lattice element formerly labeled *p2* can be removed, since it has become empty.

8 Implementation

A prototype implementation of the method was recently completed. Our tool – named KABA¹³ – is written in Java and analyzes Java Class files. This approach has the advantage that no front end is needed. Furthermore, Java is much easier to analyze than C++.

8.1 CFG and points-to analysis

The tool first reads the required class files and builds a control flow graph (CFG). Since Java byte code is stack-oriented, but our analysis needs full variable references rather than anonymous stack entries, a simple backwards analysis reconstructs the stack contents whenever necessary. If, at a certain point in the CFG, we need to know the type of, for example, the third entry from the top of stack, we explore all CFG paths backwards until three push operations have been encountered on every backward path, and collect the items pushed onto the stack by the third-last push operations. Usually, the resulting sets are unique.

For points-to analysis, we use Andersen’s method, as described in [26]. This method is quite precise, but also expensive: it has worst-case time complexity $O(n^3)$ and is very space-intensive in practice. In fact, the points-to analysis turns out to be the bottleneck of the analysis.

Points-to analysis has originally been designed for imperative languages such as C, and we had to extend it for object-oriented languages. In particular, the treatment of virtual dispatch requires special attention. The details are – for both Andersen’s and Steensgaard’s method – described in [32]. Roughly, dynamic dispatch is modeled as follows. Whenever a method call $o.m(x)$ is encountered during interprocedural iteration, the following steps are performed:

1. let $o \mapsto \{o_1, \dots, o_n\}$ and $x \mapsto \{x_1, \dots, x_m\}$ be the points-to information for o and x as obtained so far (encoded in the points-to graph). Static lookup is used to resolve any of the calls $o_1.m(x), \dots, o_n.m(x)$. Note that only “real objects” o_i must be considered; in Java these correspond to call sites of constructor functions (see 6.1).
2. Let $C_1::m(a_1), \dots, C_n::m(a_n)$ be the methods identified by static lookup. For any a_i, x_j , add edges as required by the assignments $a_i = x_j$ to the points-to graph. Furthermore, add edges as required by the implicit assignment to the **this**-pointer $C_i::m = o_i$.

¹³KABA = KlassenAnalyse mit BegriffsAnalyse [class analysis using concept analysis]. KABA is also a popular chocolate drink in Germany.

3. In case m 's return type is a class, let r_1, \dots, r_n be variables representing the return values of the $C_i::m(a_i)$ inside the method. For any assignment or similar use of the return value, such as in $y = o.m(x)$, add edges as required by the assignment $y = r_i$ to the points-to graph.
4. Continue propagation of points-to information.

8.2 Generation of table entries and implication propagation

The table is implemented as a list of bit strings, and once points-to analysis has converged, entering the member access entries into the table is straightforward. Next, the assignment and dominance implications are extracted. Extracting the dominance rules is quite expensive, because for any classes $A \leq B$ and any columns $A::m, B::m$, every row must be checked for a double entry $(x, A::m)$ and $(x, B::m)$.

Assignment implications as well as dominance implications are arranged into directed graphs. While the assignment graph may contain cycles, the dominance graph can not, as dominance edges always go from members in “lower” classes to members in “upper” classes (and the original class hierarchy of course is cycle free). Note that, while the set of assignment implications never changes, new dominance implications might be generated after applying assignment implications. Hence, the dominance graph can grow during implication propagation. This is the reason why the $O(n^2)$ method for applying implications to a table [13] cannot be used, and a fix point iteration must be used instead.

Assignment and dominance implications are applied alternatively. “Local” iteration applies assignment resp. dominance implications until they converged. “Global” iteration alternates local assignment or dominance iterations. Implication propagation proceeds in topological order, and never propagates out of a cycle until the cycle converges. Since the dominance graph is cycle free, the corresponding local iteration converges immediately.

Since our iteration is a special instance of a set-based analysis (an implication $a \rightarrow b$ corresponds to the set constraint $\sigma(a) \subseteq \sigma(b)$ resp. $\tau(a) \supseteq \tau(b)$), cycles can be collapsed as described in [11]. Application of a specific assignment implication $a \rightarrow b$ is implemented efficiently by a bitstring operation: $\sigma(b) := \sigma(b) \text{ OR } \sigma(a)$.

8.3 Interactive back end

From the final table, the lattice is computed using Ganters algorithm. Next, an off-the-shelf graph layouter is used to compute an initial layout for the lattice. The lattice is displayed by an interactive back end.

The lattice layout may be modified manually, while the system maintains lattice integrity. There are several options for the display of lattice elements labels, namely no labels at all, individual labels on request, and labels for user-defined entities only. The KABA prototype also offers some of the reengineering transformations from Section 7, namely removal of empty lattice elements, reduced lattices without pointers, highlighting

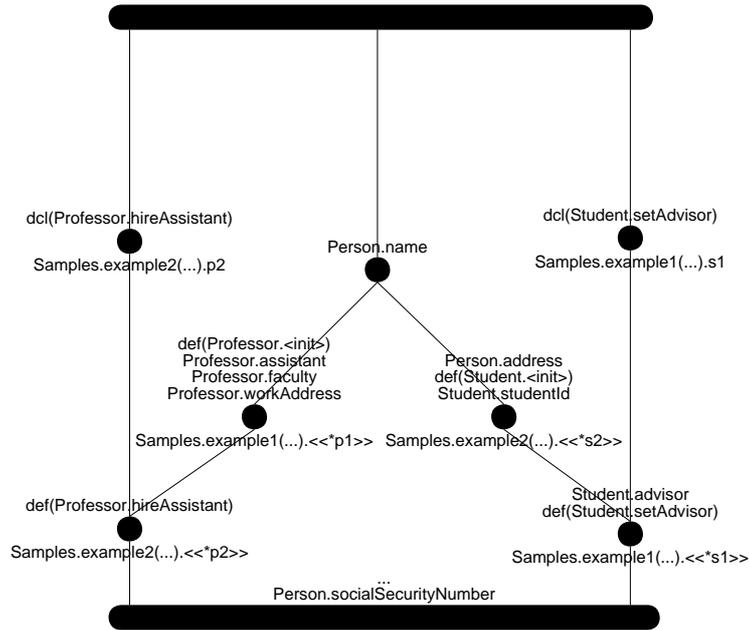


Figure 10: Java version of student/professor example

of all variables which had the same original type, and recovery of associations. Interactive lattice manipulation and code generation are not supported yet.

9 Case Studies

9.1 Students and Professors, finally

KABA was applied to several small and medium-sized Java programs. We begin with a reconsideration of the student-professor example (see Figure 1), in order to illustrate differences between C++ and Java. Figure 10 presents a screenshot. Data members appear with their fully qualified name, whereas method definitions are distinguished from method declarations. Names in “<<...>>” are names of “real” objects (heap allocation sites), and constructor methods are named “<init>”. Methods are displayed with full name, but without signature (signatures are shown only for overloaded methods). Top and bottom element are enlarged for layout reasons.

The first observation is that this lattice is different from the one in Figure 2. This is not a bug: the screenshot displays an analysis of the Java version of Figure 1. In Java, all methods are virtual, while in Figure 1, `setAdvisor` and `hireAssistant` are non-virtual. According to the table construction rules, receivers of virtual methods only need to see the method declaration, hence the definition of `Professor::hireAssistant` needs not be visible to `p2`. Consequently, data member `Professor::assistant` does not need to be visible to `p2`. Therefore a corresponding table entry is not created, and the distinction between the two rightmost lattice elements in Figure 2 disappears. The same

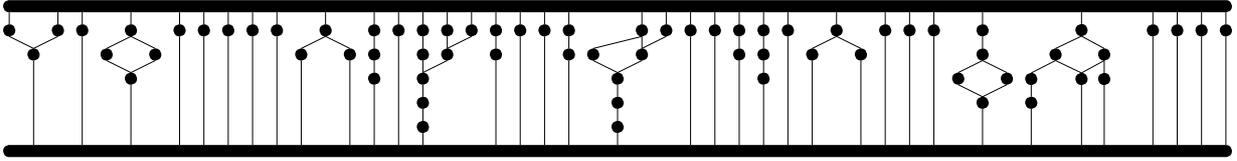


Figure 11: Lattice for graph editor program

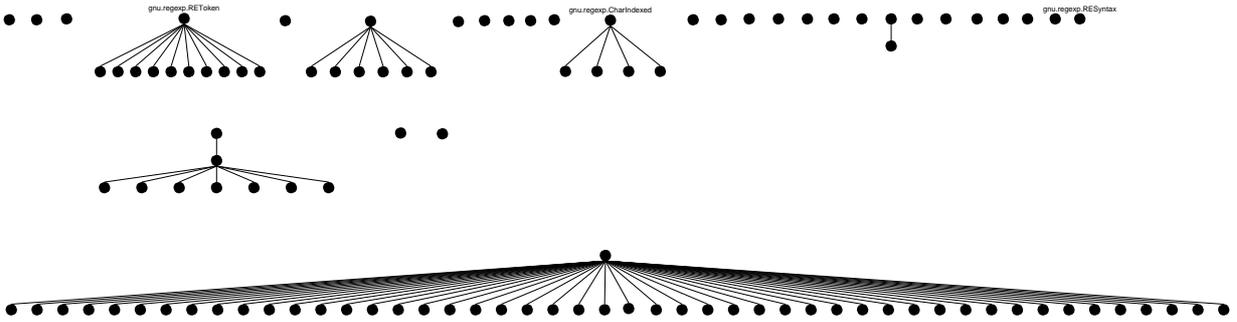


Figure 12: Original class hierarchy for “jEdit” program

argument applies to `Student::advisor` – except that the original lattice did not contain such a distinction anyway (due to the missing initialization of `Student::advisor` in the constructor, see section 7). As a result of this subtle phenomenon, the Java version of the lattice is completely symmetrical – indicating the lower semantic complexity of Java vs. C++.

9.2 An easy case

Our next example is a graph editor program (3761 LOC \ comments) with a completely flat class structure. The purpose of this experiment was to see whether KABA proposes to introduce inheritance and specialized subclasses.

The lattice (Figure 11) is horizontally decomposable into several small sublattices (actually, each sublattice corresponds to one original class). The internal structure in the sublattices comes from fine-grained pointer access patterns and should not be interpreted as an option to split classes – in particular, each substructure in the lattice (except two) has its own local bottom element, which is always an indicator that potential for introducing inheritance and splitting classes is low. Indeed, the reduced lattice without pointers replicates the original hierarchy. Thus KABA demonstrates that the original class design was good. This example shows that our approach is useful not only for reengineering, but also for ongoing quality assurance during development: it can confirm that the class design is ok.

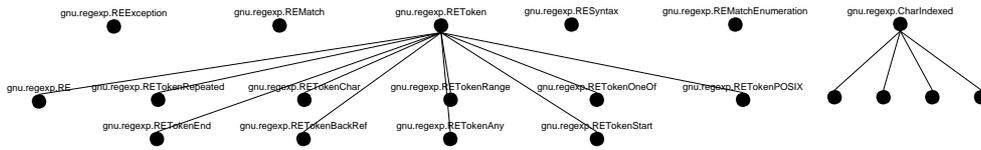


Figure 13: Subhierarchy for regular expression library

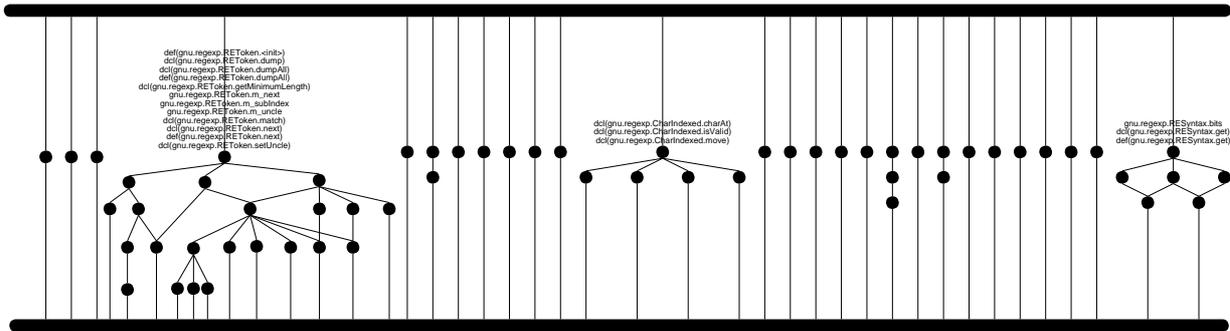


Figure 14: Lattice for “jEdit”

9.3 The GNU regular expression library

Our next example is “jEdit”, a text editor with useful features such as syntax coloring and regular expression search¹⁴. jEdit makes heavy use of a Java adaption of the GNU regular expression library, and can thus be seen as an instance of our scenario, namely that a given hierarchy is used by different applications.

Figure 12 shows the original hierarchy of all classes shipped with “jEdit”. Five separate subsystems are visible, concerned with input modes, editor commands, editor modes, regular expressions, and syntax highlighting. Several singleton classes without any inheritance relationship provide basic and auxiliary functionality. All original subhierarchies are very flat. JEdit contains more than 80 classes and almost 9000 LOC. Figure 13 shows those classes of the original hierarchy which constitute the regular expression library; these are 20 classes containing more than 3000 LOC. The superclass `REToken` has one subclass for every regular expression construct (`*`, `+`, `[]`, `$`). The programming interface class `RE` also is a subclass of `REToken`.

The reduced lattice¹⁵ computed by KABA (Figure 14) consists of several independent substructures, which correspond to the subsystems from the original hierarchy. Most of the original singleton classes, as well as the “input mode” subsystem, recur exactly in the right hand part of the lattice, showing no reengineering potential.

Most interesting however is the leftmost part of the lattice, which represents the regular

¹⁴version 1.2final, available from [http://www.gjt.org/\\$p/jedit.html](http://www.gjt.org/$p/jedit.html)

¹⁵For this and the following example we used the reduced lattice which does not show fine-grained access patterns for pointers, just for “real” objects. As explained in section 7, fine-grained pointer access patterns are not really relevant for reengineering, and the reduced lattice is still guaranteed to be operationally equivalent. The full lattices for this and the next example are about twice the size as the reduced lattices.

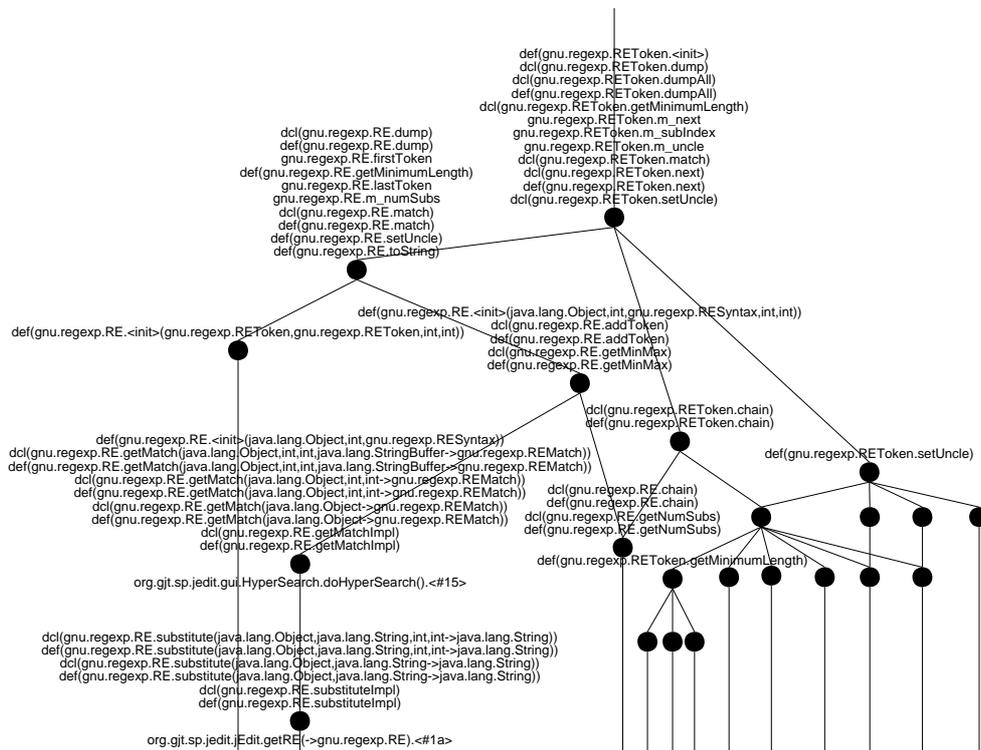


Figure 15: Details for “jEdit”: regular expression classes

expressions library (that is, all subclasses of `REToken`). It has become a complex structure (Figure 15). The right part of this structure (no labels shown) represents the classes for the different regular expression constructs; they are more fine-grained than before because the different constructs need different parts of `REToken`'s functionality. Some subclasses of the original base class `REToken` are just reproduced by KABA. However the original subclass `RE` has been distributed to 6 different nodes (left part of Figure 15). A look at the source code reveals that the class labeled `gui.HyperSearch.doHyperSearch` is the API for search without substitution, while the class below it is the API for search including substitution.

Note that the lattice displays the finest possible splittings and refactorings of classes according to possible program behaviour. For reengineering purposes, the lattice should therefore be simplified by merging lattice elements, in order to reflect software design principles. But the lattice demonstrates that the original `RE` class can be split into, say, `RE_substitution` and `RE_no_substitution`. The element labeled `gnu.regexp.RE.dump` etc. also reveals that there is a “composite” design pattern used: class `RE` stands for complex regular expression and offers method for accessing subexpressions, while the subclasses of `REToken` in the right part represent elementary regular expressions.

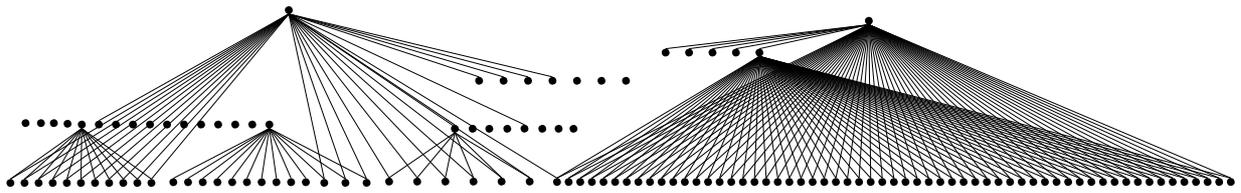


Figure 16: Original class hierarchy for the “JAS” example

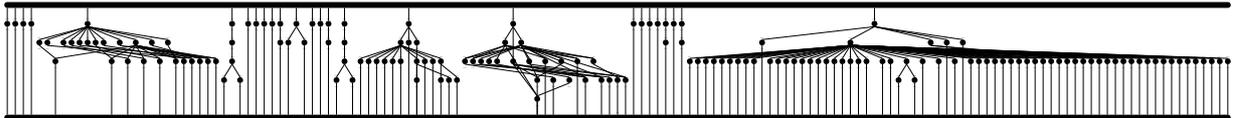


Figure 17: Lattice for “JAS”

9.4 JAS

Our next example is “JAS”, a java bytecode assembler, including a Scheme-like scripting language (about 5400 LOC)¹⁶. Its original class hierarchy is shown in Figure 16. Among various single classes and three small inheritance trees it shows a huge structure with more than 50 classes. These classes are part of the scripting language implementation. The top class is called `Obj` and all (except four) subclasses additionally implement an interface `Procedure`. Each of these subclasses represents a function like “Add” or “Sub” in the scripting language.

In the KABA lattice (Figure 17) this huge structure is reproduced basically unmodified, confirming the original design was good. But the subhierarchy with base class `Insn` (Figure 18) is interesting. All but one subclasses of `Insn` just redefine the constructor method, and these subclasses are reproduced unchanged. However the rightmost chain in figure 18, which contains all members of the original subclass `Label`, differs from the other subclasses because it does not use the methods `Insn.size` and `Insn.write`. A closer look reveals that all the other subclasses are dealing with the implementation of certain bytecode instructions, but `Label` is concerned with bytecode addressing. The implementations of `size` and `write` in `Label` are empty, so these two methods can be considered “amputated”. An even closer look reveals that the `resolve` method does not execute any useful code when called from a `Label` object. This demonstrates that the original subhierarchy should be restructured: `Label` does not share any code with the other subclasses, thus it does not need a common base class with them.

The sublattice for the subclasses of the `InsnOperand` class show a similar phenomenon (Figure 19). Two classes (`UnsignedByteWideOperand` and `IincOperand` on the left hand side) are separated from the rest, just like `Label` was. They have own implementations of the method `writePrefix`, while all other subclasses share the same implementation. A look at the source code reveals that the other subclasses use a dummy implementation of `writePrefix` which has no functionality; only the two “separated” classes on the left

¹⁶version 0.4, available from <http://www.sbktech.org/jas.html>

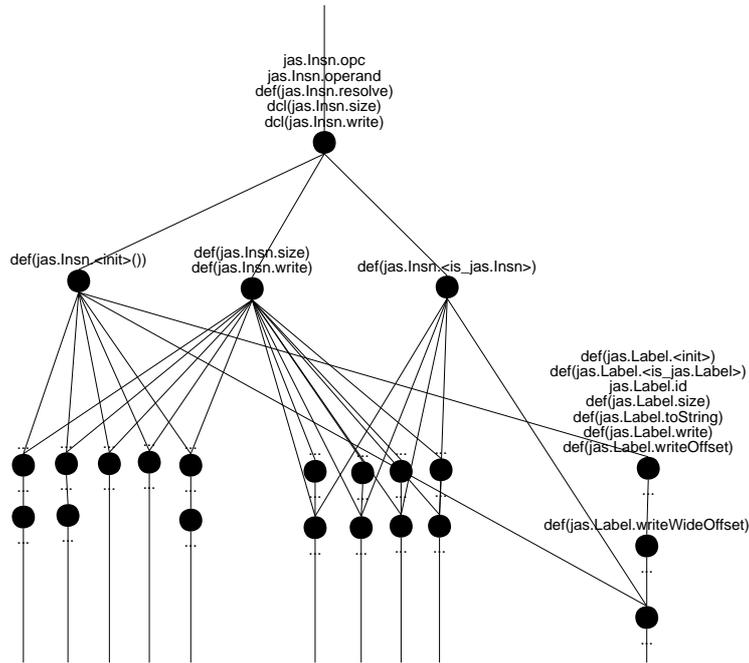


Figure 18: Details for “JAS”: substructure for `Insn`

actually have code for `writePrefix`. This demonstrates that `writePrefix` can be removed from `InsnOperand` and put into a new class, which should be the base class of the separated classes.

9.5 Additional remarks

Some additional experiments are described in the Bögemann/Streckenbach master’s thesis [6]. The preliminary experience from our experiments can be summarized as follows:

- The KABA prototype may need several hours for 20,000 LOC on a PC. We plan to switch to a native code compiler with a better garbage collector and expect to be able to analyse 50,000 LOC within reasonable time on a standard workstation.
- The bottleneck is the initial points-to analysis, which consumes up to 80% of the total analysis time. An implementation of the supposedly faster (and less precise) Steensgaard algorithm – adapted for Java – was not faster in practice, due to the necessary conservative approximation of dynamic binding.
- Most Java programs explicated a reasonable structure without high reengineering potential – probably due to the fact that these programs are quite young. In such cases, the lattice can serve as a quality metrics, demonstrating that the architecture is good.

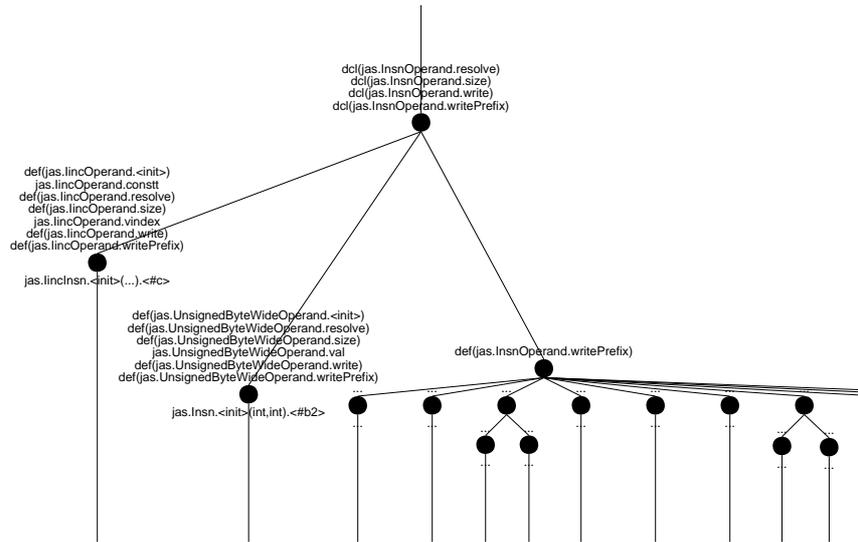


Figure 19: Details for “JAS”: substructure for `InsnOperand`

- Nevertheless, in many Java examples we found possibilities to split or refactor classes. These proposals, which are guaranteed not to alter the program behavior, would not have been possible without our unique combination of points-to analysis, type constraints, and concept lattices.
- KABA has an experimental option where member accesses from dead code will not be entered into the table. This greatly reduces the size of the lattices in many examples. From a reengineering viewpoint however, it is questionable to exclude dead code – just as it is questionable to delete dead members such as `socialSecurityNumber`.
- We did not yet exploit the structure theory of concept lattices, in particular congruences and weak congruences. (Weak) congruence classes could serve as proposals how to group classes into packages, and can be used to measure coupling and cohesion; resulting in more substantial restructuring proposals.
- The full set of reengineering transformations from section 7 is not available yet – the prototype is still incomplete.
- Of course the real “market” for the method are big old C++ programs. However the complexity of both the language and the method itself seem to prohibit an application to C++ right now. We hope that this situation will change within the next two years.

10 Related Work

10.1 Applications of concept analysis

Godin and Mili [14, 15] also use concept analysis for class hierarchy (re)design. The starting point in their approach is a set of interfaces of (collection) classes. A table is constructed that specifies for each interface the set of supported methods. The lattice derived from this table suggests how the design of a class hierarchy implementing these interfaces could be organized in a way that optimizes the distribution of methods over the hierarchy. Although Godin and Mili’s work has the same formal basis as ours, the domains under consideration are different. In [14], relations between members and classes are studied in order to improve the distribution of these members over the class hierarchy. In contrast, we study how the members of a class hierarchy are *used* in the executable code of a set of applications by examining relationships between variables and class members, and relationships among class members.

Another application of concept analysis in the domain of software engineering is the analysis of software configurations. Snelting [28] uses concept analysis to analyze systems in which the C preprocessor (CPP) is used for configuration management. The relation between code pieces and governing expressions is extracted from a source file, and the corresponding lattice visualizes interferences between configurations. Later, Lindig proved that the configuration space itself is isomorphic to the lattice of the *inverted* relation [17].

Concept analysis was also used for modularization of old software. Siff and Reps [27] investigated the relation between procedures and “features” such as usage of global variables or types. A modularization is achieved by finding elements in the lattice whose intent partitions the feature space. Lindig and Snelting [18] also analyzed the relation between procedures and global variables in legacy Fortran programs. They showed that the presence of module candidates corresponds to certain decomposition properties of the lattice (the Siff/Reps criterion being a special case).

10.2 Class hierarchy specialization and application extraction

The work in the present paper is closely related to the work on *class hierarchy specialization* by Tip and Sweeney [35, 36]. Class hierarchy specialization is a space optimization technique in which a class hierarchy and a client program are transformed in such a way that the client’s space requirements are reduced at run-time. The method of [35, 36] shares some basic “information gathering” steps with the method of the present paper¹⁷, but the subsequent steps of that method are quite different. After determining the member access and assignment operations in the program, a set of *type constraints* is computed that capture the subtype-relationships between variables and members that must be retained. These type constraints roughly correspond to the information encoded in our tables, but contrary to our current approach they correctly distinguish between multiple subobjects that have the same type. From the type constraints, a new class hierarchy is generated

¹⁷Definitions 1, 3, 4, and 7 were taken from [35, 36].

automatically. In a separate step, the resulting class hierarchy is simplified by repeatedly applying a set of simple graph rewriting rules.

In addition to the differences in the underlying algorithms, the method of [35, 36] differs from our reengineering framework in a number of ways. Class hierarchy specialization is an optimization technique that does not require any intervention by the user. In contrast, the current paper presents an *interactive* approach for analyzing the usage of a class hierarchy in order to find design problems. Reducing object size through the elimination of members is possible, but not necessarily an objective. For the purpose of restructuring it may very well be the case that an unused member should be retained in the restructured class hierarchy. The framework we presented here also allows for the analysis of a class hierarchy along with any number of programs, including none. Class hierarchy specialization customizes a class hierarchy w.r.t. a *single* client application.

Several other *application extraction* techniques for eliminating unused components from hierarchies and objects have been presented in the literature. These are primarily intended as optimizations, although they may have some value for program understanding. [2] describe an algorithm for the dynamically typed language Self that eliminates unused slots from objects (a slot corresponds to either a data member, a method, or an inheritance relation). Self is a dynamically typed language, and eliminating members from objects does not involve transforming class hierarchies.

[33] present an algorithm for slicing class hierarchies that eliminates members and inheritance relations from a C++ hierarchy. Class slicing is less powerful than specialization because it can only remove a member m from a class C if m is not used by *any* C -instance. Later, Tip and Sweeney developed Jax [34], which incorporates rapid type analysis [4] and some additional simplifications. Jax can reduce the size of class files up to 70%.

10.3 Techniques for restructuring class hierarchies

Another category of related work is that of techniques for restructuring class hierarchies for the sake of improving design, improving code reuse, and enabling reuse. The overview article [7] presents 18 different methods, many of them process-centered or dynamic analyses. The probably most well-known method for static restructuring was introduced by Opdyke and Johnson [21, 20]. They present a number of behavior-preserving transformations on class hierarchies, which they refer to as *refactorings*. The goal of refactoring is to improve design and enable reuse by “factoring out” common abstractions. This involves steps such as the creation of new superclasses, moving around methods and classes in a hierarchy, and a number of similar steps. Our techniques for analyzing the usage of a class hierarchy to find design problems is in our opinion complimentary to the techniques of [21, 20].

Moore [19] presents a tool that automatically restructures inheritance hierarchies and refactors methods in Self programs. The goal of this restructuring is to maximize the sharing of expressions between methods, and the sharing of methods between objects in order to obtain smaller programs with improved code reuse. Since Moore is studying a dynamically typed language without explicit class definitions, a number of complex issues related to preserving the appropriate subtype-relationships between types of variables do

not arise in his setting.

An interesting approach is that of Astudillo [3]. He argues that from an evolutionary viewpoint, subclasses may not only add or redefine members, but also lose or “amputate” members. This approach violates fundamental type-theoretic properties of object-oriented programming, but has the advantage that well-known algorithms for the reconstruction of biological taxonomies can be used. Astudillo argues that his hierarchies are more “natural” than those which stick to the principles of type conformance and contravariance.

11 Conclusions and Future Work

We have presented a method for finding design problems in a class hierarchy by analyzing the *usage* of the hierarchy by a set of applications. This method constructs a *concept lattice* in which relationships between variables and class members are made explicit, and where information that members and variables have in common is “factored out”. We have shown the technique to be capable of finding design anomalies such as class members that are redundant or that can be moved into a derived class. In addition, situations where it is appropriate to split a class can be detected. We have suggested how these techniques can be incorporated into interactive tools for maintaining and restructuring class hierarchies.

Our analysis is perhaps the most expensive analysis of object-oriented programs available at the moment. But it is also one of the most powerful methods, due to its unique combination of points-to analysis, type constraints, and concept lattices. The method includes classic analyses such as dead members or useless variables as special cases. Our preliminary case studies have indicated the usefulness of the analysis as a basis for reengineering, but the method can also be used for quality assessment during initial development. It turned out that the Java examples we analysed were all reasonably well structured, but that the method nevertheless discovered many possibilities for refactoring – while at the same time guaranteeing that program behavior is unchanged.

The present article has focused on foundational aspects and preliminary case studies. Future work will concentrate on questions of scale-up (apply the method to large programs), better interactive support for restructuring (in particular, moving members up, variables down, merge lattice elements, and integrate background knowledge), utilizing the structure theory of concept lattices, and, eventually, develop a version for C++. Of course, big old C++ programs are the real market for the method. It remains to be seen whether an efficient implementation for full C++ can be achieved.

Acknowledgements. Andreas Bögemann and Mirko Streckenbach did a great job with the prototype implementation [6], and Mirko’s help with the experiments and with improving the prototype was indispensable. This work is supported by Deutsche Forschungsgemeinschaft, grant Sn11/7-1.

References

- [1] Information Processing Systems Accredited Standards Committee X3. Working paper for draft proposed international standard for information systems—programming language C++. Doc. No. X3J16/97-0108. Draft of 25 November 1997.
- [2] Ole Agesen and David Ungar. Sifting out the gold: Delivering compact applications from an exploratory object-oriented programming environment. In *Proceedings of the Ninth Annual Conference on Object-Oriented Programming Systems, Languages, and Applications (OOPSLA '94)*, pages 355–370, Portland, OR, 1994. *ACM SIGPLAN Notices* 29(10).
- [3] H. Astudillo. Maximizing object reuse with a biological metaphor. *Theory and practice of object systems*, 3(4):235–251, 1997.
- [4] David F. Bacon and Peter F. Sweeney. Fast static analysis of C++ virtual function calls. In *Proceedings of the Eleventh Annual Conference on Object-Oriented Programming Systems, Languages, and Applications (OOPSLA '96)*, pages 324–341, San Jose, CA, 1996. *ACM SIGPLAN Notices* 31(10).
- [5] Gary Birkhoff. *Lattice Theory*. American Mathematical Society, 1940.
- [6] Andreas Bögemann and Mirko Streckenbach. KABA: Reengineering class hierarchies using concept lattices. Master's thesis, Technische Universität Braunschweig, Germany, May 1999.
- [7] Eduardo Casais. Reengineering of object-oriented legacy systems. *Journal of object-oriented programming*, pages 45–52, January 1998.
- [8] J.-D. Choi, M. Burke, and P. Carini. Efficient flow-sensitive interprocedural computation of pointer-induced aliases and side effects. In *Conference Record of the Twentieth ACM Symposium on Principles of Programming Languages*, pages 232–245. ACM, 1993.
- [9] J.D. Choi, D. Grove, M. Hind, and V. Sarkar. Efficient and precise modelling of exceptions for the analysis of java programs. In *Proc. PASTE '99*. ACM, 1999. to appear.
- [10] B. Davey and H. Priestley. *Introduction to lattices and order*. Cambridge University Press, 1990.
- [11] Manuel Fähndrich, Jeffrey Foster, Zhendong Su, and Alexander Aiken. Partial online cycle elimination in inclusion constraint graphs. In *Proceedings of the ACM SIGPLAN'98 Conference on Programming Language Design and Implementation*, pages 85–96, Montreal, Canada, June 1998. *ACM SIGPLAN Notices* 33(6).

- [12] Petra Funk, Anke Lewien, and Gregor Snelting. Algorithms for concept lattice decomposition and their applications. Technical Report 95-05, TU Braunschweig, FB Informatik, 1998.
- [13] Bernhard Ganter and Rudolf Wille. *Formal Concept Analysis - Mathematical Foundations*. Springer Verlag, 1999.
- [14] Robert Godin and Hamed Mili. Building and maintaining analysis-level class hierarchies using galois lattices. In *Proceedings of the Eighth Annual Conference on Object-Oriented Programming Systems, Languages, and Applications (OOPSLA '93)*, pages 394–410, Washington, DC, 1993. *ACM SIGPLAN Notices* 28(10).
- [15] Robert Godin, Hamed Mili, Guy W. Mineau, Rokia Missaoui, Amina Arfi, and Thuy-Tien Chau. Design of class hierarchies based on concept (galois) lattices. *Theory and Practice of Object Systems*, 4(2):117–134, 1998.
- [16] Maren Krone and Gregor Snelting. On the inference of configuration structures from source code. In *Proceedings of the 1994 International Conference on Software Engineering (ICSE'94)*, pages 49–57, Sorrento, Italy, May 1994.
- [17] Christian Lindig. Analyse von Softwarevarianten. Technical Report 98-02, TU Braunschweig, FB Informatik, 1998.
- [18] Christian Lindig and Gregor Snelting. Assessing modular structure of legacy code based on mathematical concept analysis. In *Proceedings of the 1997 International Conference on Software Engineering (ICSE'97)*, pages 349–359, Boston, MA, May 1997.
- [19] Ivan Moore. Automatic inheritance hierarchy restructuring and method refactoring. In *Proceedings of the Eleventh Annual Conference on Object-Oriented Programming Systems, Languages, and Applications (OOPSLA '96)*, pages 235–250, San Jose, CA, 1996. *ACM SIGPLAN Notices* 31(10).
- [20] W. Opdyke and R. Johnson. Creating abstract superclasses by refactoring. In *ACM 1993 Computer Science Conference*, 1993.
- [21] William F. Opdyke. *Refactoring Object-Oriented Frameworks*. PhD thesis, University Of Illinois at Urbana-Champaign, 1992.
- [22] Hemant D. Pande and Barbara G. Ryder. Data-flow-based virtual function resolution. In *Proceedings of the Third International Symposium on Static Analysis (SAS'96)*, pages 238–254, September 1996. Springer-Verlag LNCS 1145.
- [23] William Pugh and David Wonnacot. Constraint-based array dependence analysis. *ACM Transactions on Programming Languages and Systems*, 20(3):635–678, May 1998.

- [24] G. Ramalingam and H. Srinivasan. A member lookup algorithm for C++. In *Proceedings of the ACM SIGPLAN'97 Conference on Programming Language Design and Implementation*, pages 18–30, Las Vegas, NV, 1997.
- [25] Jonathan G., Jr. Rossie and Daniel P. Friedman. An algebraic semantics of subobjects. In *Proceedings of the Tenth Annual Conference on Object-Oriented Programming Systems, Languages, and Applications (OOPSLA '95)*, pages 187–199, Austin, TX, 1995. *ACM SIGPLAN Notices* 30(10).
- [26] Marc Shapiro and Susan Horwitz. Fast and accurate flow-insensitive points-to analysis. In *Conference Record of the Twenty-Fourth ACM Symposium on Principles of Programming Languages*, pages 1–14, Paris, France, 1997.
- [27] M. Siff and T. Reps. Identifying modules via concept analysis. In *Proc. International Conference on Software Maintenance*, pages 170–179, Bari, Italy, 1997.
- [28] Gregor Snelting. Reengineering of configurations based on mathematical concept analysis. *ACM Transactions on Software Engineering and Methodology*, 5(2):146–189, April 1996.
- [29] Gregor Snelting. Concept analysis – a new framework for program understanding. In *Proc. ACM SIGPLAN/SIGSOFT Workshop on Program Analysis for Software Tools and Engineering (PASTE)*, pages 1–10, Montreal, Canada, June 1998. *ACM SIGPLAN Notices* 33(7).
- [30] Gregor Snelting and Frank Tip. Reengineering class hierarchies using concept analysis. In *Proc. ACM SIGSOFT Symposium on the Foundations of Software Engineering*, pages 99–110, Orlando, FL, November 1998.
- [31] Bjarne Steensgaard. Points-to analysis in almost linear time. In *Proceedings of the Twenty-Third ACM Symposium on Principles of Programming Languages*, pages 32–41, St. Petersburg, FL, January 1996.
- [32] Mirko Streckenbach and Gregor Snelting. Points-to analysis for object-oriented languages. Technical report, Universität Passau, Fakultät für Informatik, 1999. In preparation.
- [33] Frank Tip, Jong-Deok Choi, John Field, and G. Ramalingam. Slicing class hierarchies in C++. In *Proceedings of the Eleventh Annual Conference on Object-Oriented Programming Systems, Languages, and Applications (OOPSLA '96)*, pages 179–197, San Jose, CA, 1996. *ACM SIGPLAN Notices* 31(10).
- [34] Frank Tip, Chris Laffra, Peter F. Sweeney, and David Streeter. Practical experience with an application extractor for java. In *Proc. OOPSLA '99*, Denver, November 1999. to appear.

- [35] Frank Tip and Peter Sweeney. Class hierarchy specialization. In *Proceedings of the Twelfth Annual Conference on Object-Oriented Programming Systems, Languages, and Applications (OOPSLA '97)*, pages 271–285, Atlanta, GA, 1997. *ACM SIGPLAN Notices* 32(10).
- [36] Frank Tip and Peter F. Sweeney. Class hierarchy specialization. Technical Report RC21111, IBM T.J. Watson Research Center, February 1998. Submitted for publication.
- [37] Rudolf Wille. Restructuring lattice theory: an approach based on hierarchies of concepts. *Ordered Sets*, pages 445–470, 1982.

ppendix

This appendix demonstrates that a method and its **this**-pointer will always appear together in the lattice, and that arbitrary pointers appear above any object they point to.

Lemma. For any $a \in \mathcal{A}$, $\mu(a) = \bigvee_{(o, a) \in T} \gamma(o)$ **Proof.** This is true for every concept lattice by construction [13].

Lemma. For any $\text{def}(C::f) \in \text{MemberDefs}(\mathcal{P})$, we have that:

$$\gamma(*C::f) \geq \mu(\text{def}(C::f))$$

Proof.

By the preceding lemma,

$$\mu(\text{def}(C::f)) = \bigvee_{(x, \text{def}(C::f)) \in T} \gamma(x)$$

Furthermore, for any x such that $(x, \text{def}(C::f)) \in T$, there must be a method call $x.f()$ (otherwise the table entry would not exist). This method call causes an implicit assignment to f 's **this** pointer, generating an assignment dominance $*C::f \rightarrow x$, which enforces $\gamma(*C::f) \geq \gamma(x)$. Since this is true for all x calling f , it is also true for their supremum:

$$\gamma(*C::f) \geq \bigvee_{(x, \text{def}(C::f)) \in T} \gamma(x)$$

Combining both statements, we obtain $\gamma(*C::f) \geq \mu(\text{def}(C::f))$.

The last lemma shows that a method always appears below its **this**-pointer, and without the **this**-rule, they will indeed appear at different elements in the lattice if method $C::f$ does not access itself (i.e., is non-recursive).

The **this**-rule enforces $\gamma(*C::f) \leq \mu(\text{def}(C::f))$, and together with the lemma we may conclude the following.

Proposition. For any $\text{def}(C::m) \in \text{MemberDefs}(\mathcal{P})$ we have that:

$$\gamma(*C::f) = \mu(\text{def}(C::f))$$

Hence methods and their **this**-pointers appear together in the lattice. For pointers in general, only a weaker result can be established: any pointer always appears above any object it may point to.

Proposition.

$$\langle p, v \rangle \in \text{PointsTo}(\mathcal{P}) \implies \gamma(p) \geq \gamma(v)$$

Proof. (In this proof, we use member m also as a shorthand for $\text{def}(X::m)$ or $\text{dcl}(X::m)$.)
By the basic properties of concept lattices, it is enough to show

$$\forall m \in \text{MemberDcls}(\mathcal{P}) : (p, m) \in T \implies (v, m) \in T$$

because this implication will force $\gamma(p) \geq \gamma(v)$. So let $(p, m) \in T$. By Definition 5, this implies $\langle m, p \rangle \in \text{MemberAccess}(\mathcal{P})$ and therefore $p.m()$ must occur in \mathcal{P} . Since $\langle p, v \rangle \in \text{PointsTo}(\mathcal{P})$, by definition 4 (case 3), $\langle m, v \rangle \in \text{MemberAccess}(\mathcal{P})$ and therefore $(v, m) \in T$.