

A Relational Static Semantics for Call Graph Construction

Xilong Zhuo and Chenyi Zhang

Jinan University

Abstract. The problem of resolving virtual method and interface calls in object-oriented languages has been a long standing challenge to the program analysis community. The complexities are brought by various reasons, such as increased levels of class inheritance and polymorphism. In this paper, we propose a new approach called type flow analysis that represent the propagation of type information between program variables by a group of relations without the help of a heap abstraction. We prove that regarding the precision on reachability of class information to a variable, our method produces results equivalent to that one can derive from a points-to analysis. Moreover, in practice, our method consumes lower time and space usage, as supported by the experimental results.

1 Introduction

For object-oriented programming languages, virtual methods (or functions) are those declared in a base class but are meant to be overridden in different child classes. Statically determine a set of methods that may be invoked at a call site is important to program optimization, from result of which a subsequent optimization may reduce the cost of virtual function calls or perform method inlining if target method forms a singleton set, and one may also remove methods that are never called by any call sites, or produce a call graph which can be useful in other optimization processes.

Efficient solutions, such as Class Hierarchy Analysis (CHA) [?,?], Rapid Type Analysis (RTA) [?] and Variable Type Analysis (VTA) [?], conservatively assign each variable a set of class definitions, with relatively low precision. Alternatively, with the help of an abstract heap, one may take advantage of points-to analysis [?] to compute a set of object abstractions that a variable may refer to, and resolve the receiver classes in order to find associated methods at call sites.

The algorithms used by CHA, RTA and VTA are conservative, which aim to provide an efficient way to resolve calling edges, and which usually runs linear-time in the size of a program, by focusing on the types that are collected at the receiver of a call site. For instance, let x be a variable of declared class A , then at a call site $x.m()$, CHA will draw a call edge from this call site to method $m()$ of class A and every definition $m()$ of a class that extends A . In case class A does not define $m()$, a call edge to an ancestor class that defines $m()$ will also be included. For a variable x of declared interface I , CHA will draw a call

```

class A{
    A f;
    void m(){
        return this.f;
    }
}
class B extends A{}
class C extends A{}
1:  A x = new A(); //O_1
2:  B b = new B(); //O_2
3:  A y = new A(); //O_3
4:  C c = new C(); //O_4
5:  x.f = b;
6:  y.f = c;
7:  z = x.m();

```

Fig. 1. An example that compares precision on type flow in a program.

| Statement | VTA fact |
|---------------------------|---|
| $A \ x = \text{new } A()$ | $x \leftarrow A$ |
| $B \ b = \text{new } B()$ | $b \leftarrow B$ |
| $A \ y = \text{new } A()$ | $y \leftarrow A$ |
| $C \ c = \text{new } C()$ | $c \leftarrow C$ |
| $x.f = b$ | $A.f \leftarrow b$ |
| $y.f = c$ | $A.f \leftarrow c$ |
| $z = x.m()$ | $A.m.this \leftarrow x$ $A.m.return \leftarrow A.f$ $z \leftarrow A.m.return$ |

Fig. 2. VTA facts on the example

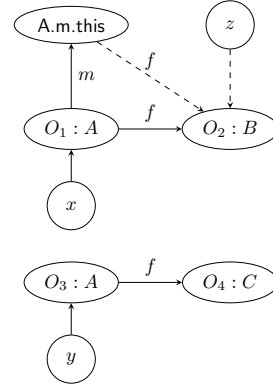


Fig. 3. Points-to results on the example

edge from this call site to every method of name $m()$ defined in class X that implements I . We write $CHA(x, m)$ for the set of methods that are connected from call site $x.m()$ as resolved by Class Hierarchy Analysis (CHA). Rapid Type Analysis (RTA) is an improvement from CHA which resolves call site $x.m()$ to $CHA(x, m) \cap inst(P)$, where $inst(P)$ stands for the set of methods of classes that are instantiated in the program.

Variable Type Analysis (VTA) is a further improvement. VTA defines a node for each variable, method, method parameter and field. Class names are treated as values and propagation of such values between variables work in the way of value flow. As shown in Figure 2, the statements on line 1 – 4 initialize type information for variables x , y , b and c , and statements on line 5 – 7 establish value flow relations. Since both x and y are assigned type A , $x.f$ and $y.f$ are both represented by node $A.f$, thus the set of types reaching $A.f$ is now $\{B, C\}$. (Note this is a more precise result than CHA and RTA which assigns $A.f$ with the set $\{A, B, C\}$.) Since $A.m.this$ refers to x , $this.f$ inside method $A.m()$ now refers to $A.f$. Therefore, through $A.m.return$, z receives $\{B, C\}$ as its final set of reaching types.

The result of a context-insensitive subset based points-to analysis [?] creates a heap abstraction of four objects (shown on line 1 – 4 of Figure 1 as well as in the ellipses in Figure 3). These abstract objects are then inter-connected via

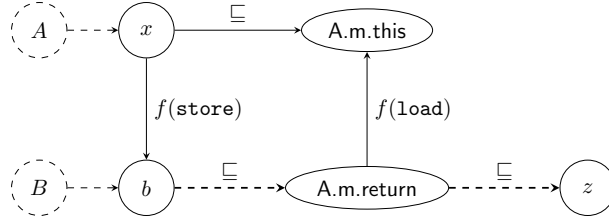


Fig. 4. Type Flow Analysis for variable z in the example

field store access defined on line 5 – 6. The derived field access from $A.m.this$ to O_2 is shown in dashed arrow. By return of the method call $z = x.m()$, variable z receives O_2 of type B from $A.m.this.f$, which gives a more precise set of reaching type for variable z .

From this example, one may conclude that the imprecision of VTA in comparison with points-to analysis is due to the over abstraction of object types, such that O_1 and O_3 , both of type A , are treated under the same type. Nevertheless, points-to analysis requires to construct a heap abstraction, which brings in extra information, especially when we are only interested in the set of reaching types of a variable.

In this paper we introduce a relational static semantics called Type Flow Analysis (TFA) on program variables and field accesses. Different from VTA, besides a binary value flow relation \sqsubseteq on the variable domain \mathbf{VAR} , where $x \sqsubseteq y$ denotes all types that flow to x also flow to y , we also build a ternary field store relation $\rightarrow \subseteq \mathbf{VAR} \times \mathcal{F} \times \mathbf{VAR}$ to trace the *load* and *store* relationship between variables via field accesses. This provides us additional ways to enlarge the relations \sqsubseteq as well as \rightarrow . Taking the example from Figure 1, we are able to collect the store relation $x \xrightarrow{f} b$ from line 5. Since $x \sqsubseteq A.m.this$, together with the implicit assignment which loads f of $A.m.return$, we further derives $b \sqsubseteq A.m.return$ and $b \sqsubseteq z$ (dashed thick arrows in Figure 4). Therefore, we assign type B to variable z . The complete reasoning pattern is depicted in Figure 4. Nevertheless, one cannot derive $c \sqsubseteq z$ in the same way.

We have proved that in the context-insensitive inter-procedural setting, TFA is as precise as the subset based points-to analysis regarding type related information. Since points-to analysis can be enhanced with various types of context-sensitivity on variables and objects (e.g., call-site-sensitivity $[?, ?]$, object-sensitivity $[?, ?, ?]$ and type-sensitivity $[?]$), our type flow analysis will only require to consider contexts on variables, which is left for future work. The context-insensitive type flow analysis has been implemented in the Soot framework $[?]$, and the implementation has been tested on a collection of benchmark programs from SPECjvm2008 $[?]$ and DaCapo $[?]$. The initial experimental result has shown that TFA has similar or less runtime than CHA $[?]$, but has precision comparable to that of a points-to analysis.

2 Type Flow Analysis

We define a core calculus consisting of most of the key object-oriented language features, as shown in Figure 5. A program is defined as a code base \overline{C} (i.e., a collection of class definitions) with statement s to be evaluated. To run a program, one may assume that s is the default (static) entry method with local variable declarations \overline{D} , similar to e.g., Java and C++, which may differ in specific language designs. We define a few auxiliary functions. $fields$ maps class names to its fields, $methods$ maps class names to its defined or inherited methods, and $type$ provides types (or class names) for objects. Given class c , if $f \in fields(c)$, then $f_{type}(c, f)$ is the defined class type of field f in c . Similarly, give an object o , if $f \in fields(type(o))$, then $o.f$ may refer to an object of type $f_{type}(type(o), f)$, or an object of any of its subclass at runtime. Write \mathcal{C} for the set of classes, \mathbf{OBJ} for the set of objects, \mathcal{F} for the set of fields and \mathbf{VAR} for the set of variables that appear in a program.

$$\begin{aligned}
C &::= \text{class } c [\text{extends } c] \{ \overline{F}; \overline{M} \} \\
F &::= c \ f \\
D &::= c \ z \\
M &::= m(x) \{ \overline{D}; s; \text{return } x' \} \\
s &::= e \mid x = \text{new } c \mid x = e \mid x.f = y \mid s; s \\
e &::= \text{null} \mid x \mid x.f \mid x.m(y) \\
prog &::= \overline{C}; \overline{D}; s
\end{aligned}$$

Fig. 5. Abstract syntax for the core language.

In this simple language we do not model common types (e.g., `int` and `float`) that are irrelevant to our analysis, and we focus on the reference types which form a class hierarchical structure. Similar to Variable Type Analysis (VTA), we assume a context insensitive setting, such that every variable can be uniquely determined by its name together with its enclosing class and methods. For example, if a local variable x is defined in method m of class c , then $c.m.x$ is the unique representation of that variable. Therefore, it is safe to drop the enclosing class and method name if it is clear from the context. In general, we have the following types of variables in our analysis: (1) local variables, (2) method parameters, (3) this reference of each method, all of which are syntactically bounded by their enclosing methods and classes.

We enrich the variable type analysis with the new type flow analysis by using three relations, a partial order on variables $\sqsubseteq \subseteq \mathbf{VAR} \times \mathbf{VAR}$, a type flow relation $--\rightarrow \subseteq \mathcal{C} \times \mathbf{VAR}$, as well as a field access relation $\longrightarrow \subseteq \mathbf{VAR} \times \mathcal{F} \times \mathbf{VAR}$, which are initially given as follows.

Definition 1. (*Base Relations*) We have the following base facts for the three relations.

1. $c --\rightarrow x$ if there is a statement $x = \text{new } c$;

2. $y \sqsubseteq x$ if there is a statement $x = y$;
3. $x \xrightarrow{f} y$ if there is a statement $x.f = y$.

Intuitively, $c \dashrightarrow x$ means variable x may have type c (i.e., c flows to x), $y \sqsubseteq x$ means all types flow to y also flow to x , and $x \xrightarrow{f} y$ means from variable x and field f one may access variable y .¹ These three relations are then extended by the following rules.

Definition 2. (*Extended Relations*)

1. For all statements $x = y.f$, if $y \xrightarrow{f}^* z$, then $z \sqsubseteq^* x$.
2. $c \dashrightarrow^* y$ if $c \dashrightarrow y$, or $\exists x \in \text{VAR} : c \dashrightarrow^* x \wedge x \sqsubseteq^* y$;
3. $y \sqsubseteq^* x$ if $x = y$ or $y \sqsubseteq x$ or $\exists z \in \text{VAR} : y \sqsubseteq^* z \wedge z \sqsubseteq^* x$;
4. $y \xrightarrow{f}^* z$ if $\exists x \in \text{VAR} : x \xrightarrow{f} z \wedge (\exists z' \in \text{VAR} : z' \sqsubseteq^* y \wedge z' \sqsubseteq^* x)$;
5. The type information is used to resolve each method call $x = y.m(z)$.

$$\forall c \dashrightarrow^* y : \quad m(z')\{\dots \text{return } x'\} \in \text{methods}(c) : \begin{cases} z \sqsubseteq^* c.m.z' \\ c \dashrightarrow^* c.m.\text{this} \\ c.m.x' \sqsubseteq^* x \end{cases}$$

The final relations are the least relations that satisfy constraints of Definition 2. Comparing to VTA [?], we do not have field reference $c.f$ for each class c defined in a program. Instead, we define a relation that connects the two variable names and one field name. (to compare in the motivating example) Although the three relations are inter-dependent, one may find that without method call (i.e., Definition 2.4), one may identify that the two relations \dashrightarrow^* (field access) and \sqsubseteq^* (variable partial order) can be determined without considering the type flow relation \dashrightarrow^* .

In order to compare the precision of TFA with points-to analysis, we present a brief list of the classic subset-based points-to rules for our language in Figure 6, where $\text{param}(\text{type}(o), m)$, $\text{this}(\text{type}(o), m)$ and $\text{return}(\text{type}(o), m)$ refer to the formal parameter, this reference and return variable of the method m of the class for which object o is declared, respectively.

To this end we present the first result of the paper, which basically says type flow analysis has the same precision regarding type based check, such as call site resolution and cast failure check, when comparing with the points-to analysis.

Theorem 1 *In a context-insensitive analysis, for all variables x and classes c , $c \dashrightarrow^* x$ iff there exists an object abstraction o of c such that $o \in \Omega(x)$.*

Proof. (sketch) For a proof sketch, first we assume every object creation site $x = \text{new } c_i$ at line i defines a mini-type c_i , and if the theorem is satisfied in this setting, a subsequent merging of mini-types into classes will preserve the result.

Moreover, we only need to prove the intraprocedural setting which is the result of Lemma 1. Because if in the intraprocedural setting the two systems

¹ Note that VTA treats statement $x.f = y$ as follows. For each class c that flows to x which defines field f , VTA assigns all types that flow to y also to $c.f$.

| statement | Points-to constraints |
|---------------------|--|
| $x = \text{new } c$ | $o_i \in \Omega(x)$ |
| $x = y$ | $\Omega(y) \subseteq \Omega(x)$ |
| $x = y.f$ | $\forall o \in \Omega(y) : \Phi(o, f) \subseteq \Omega(x)$ |
| $x.f = y$ | $\forall o \in \Omega(x) : \Omega(y) \subseteq \Phi(o, f)$ |
| $x = y.m(z)$ | $\forall o \in \Omega(y) : \begin{cases} \Omega(z) \subseteq \Omega(\text{param}(\text{type}(o), m)) \\ \Omega(\text{this}(\text{type}(o), m)) = \{o\} \\ \forall x' \in \text{return}(\text{type}(o), m) : \\ \Omega(x') \subseteq \Omega(x) \end{cases}$ |

Fig. 6. Constraints for points-to analysis.

have the same smallest model for all methods, then at each call site $x = y.m(a)$ both analyses will assign y the same set of classes and thus resolve the call site to the same set of method definitions, and as a consequence, each method body will be given the same set of extra conditions, thus all methods will have the same initial condition for the next round iteration. Therefore, both inter-procedural systems will eventually stabilize at the same model. \square

Lemma 1. *In a context-insensitive intraprocedural analysis where each class c only syntactically appears once in the form of $\text{new } c$, for all variables x and classes c , $c \dashrightarrow^* x$ iff there exists an object abstraction o of type c such that $o \in \Omega(x)$.*

Proof. Since the points-to constraints define the smallest model (Ω, Φ) with $\Omega : \text{VAR} \rightarrow \text{OBJ}$ and $\Phi : \text{OBJ} \times \mathcal{F} \rightarrow \mathcal{P}(\text{OBJ})$, and the three relations of type flow analysis also define the smallest model that satisfies Definition 1 and Definition 2, we prove that every model of points-to constraints is also a model of TFA, and vice versa. Then the least model of both systems must be the same, as otherwise it would lead to contradiction.

(\star) For the ‘only if’ part (\Rightarrow), we define $\text{Reaches}(x) = \{c \mid c \dashrightarrow^* x\}$, and assume a bijection $\xi : \mathcal{C} \rightarrow \text{OBJ}$ that maps each class c to the unique object o that is defined (and $\text{type}(o) = c$). Then we construct a function $\text{Access} : \mathcal{C} \times \mathcal{F} \rightarrow \mathcal{P}(\mathcal{C})$ and show that $(\xi(\text{Reaches}), \xi(\text{Access}))$ satisfies the points-to constraints. Define $\text{Access}(c, f) = \{c' \mid x \xrightarrow{f}^* y \wedge c \in \text{Reaches}(x) \wedge c' \in \text{Reaches}(y)\}$. We prove the following cases according to the top four points-to constraints in Figure 6.

- For each statement $x = \text{new } c$, we have $\xi(c) \in \xi(\text{Reaches}(x))$;
- For each statement $x = y$, we have $\text{Reaches}(y) \subseteq \text{Reaches}(x)$ and $\xi(\text{Reaches}(y)) \subseteq \xi(\text{Reaches}(x))$;
- For each statement $x.f = y$, we have $x \xrightarrow{f} y$, then by definition for all $c \in \text{Reaches}(x)$, and $c' \in \text{Reaches}(y)$, we have $c' \in \text{Access}(c, f)$, therefore $\xi(c') \in \xi(\text{Reaches}(y))$ we have $\xi(c') \in \xi(\text{Access}(\xi(c), f))$.
- For each statement $x = y.f$, let $c \in \text{Reaches}(y)$, we need to show $\xi(\text{Access}(c, f)) \subseteq \xi(\text{Reaches}(x))$, or equivalently, $\text{Access}(c, f) \subseteq \text{Access}(x)$. Let $c' \in \text{Access}(c, f)$, then by definition, there exist z, z' such that $c \in \text{Reaches}(z)$, $c' \in \text{Reaches}(z')$

and $z \xrightarrow{f^*} z'$. By $c \in \text{Reaches}(y)$ and Definition 2.4, we have $y \xrightarrow{f^*} z'$. Then by Definition 2.1, $z' \sqsubseteq^* x$. Therefore $c' \in \text{Reaches}(x)$.

(\star) For the ‘if’ part (\Leftarrow), let (Ω, Φ) be a model that satisfies all the top four constraints defined in Figure 6, and a bijection $\xi : \mathcal{C} \rightarrow \mathbf{OBJ}$, we show the following constructed relations satisfy value points-to.

- For all types c and variables x , $c \dashrightarrow^* x$ if $\xi(c) \in \Omega(x)$;
- For all variables x and y , $x \sqsubseteq^* y$ if $\Omega(x) \subseteq \Omega(y)$;
- For all variables x and y , and for all fields f , $x \xrightarrow{f^*} y$ if for all $o_1, o_2 \in \mathbf{OBJ}$ such that $o_1 \in \Omega(x)$ and $o_2 \in \Omega(y)$ then $o_2 \in \Phi(o_1, f)$.

We check the following cases for the three relations \dashrightarrow^* , \sqsubseteq^* and \rightarrow defined from points-to.

- For each statement $x = \text{new } c$, we have $\xi(c) \in \Omega(x)$, so $c \dashrightarrow^* x$ by definition.
- For each statement $x = y$, we have $\Omega(y) \subseteq \Omega(x)$, therefore $y \sqsubseteq^* x$ by definition.
- For each statement $x.f = y$, we have for all $o_1 \in \Omega(x)$ and $o_2 \in \Omega(y)$, $o_2 \in \Phi(o_1, f)$, which derives $x \xrightarrow{f^*} y$ by definition.
- For each statement $x = y.f$, given $y \xrightarrow{f^*} z$, we need to show $z \sqsubseteq^* x$. Equivalently, by definition we have for all $o_1 \in \Omega(y)$ and $o_2 \in \Omega(z)$, $o_2 \in \Phi(o_1, f)$. Since points-to relation gives $\Phi(o_1, f) \subseteq \Omega(x)$, we have $o_2 \in \Omega(x)$, which derives $\Omega(z) \subseteq \Omega(x)$, the definition of $z \sqsubseteq^* x$.
- The proof for the properties in the rest of Definition 2 are related to transitivity of the three TFA relations, which are straightforward. We leave them for interested readers. \square

3 Implementation and Optimization

The analysis algorithms are in Java, which we implement in the Soot framework [?], the most popular static analysis framework for Java. The three base relations (i.e., \dashrightarrow^* , \sqsubseteq^* and \rightarrow) of Definition 1 are extracted from Soot’s intermediate representation and the extended relations (i.e., \dashrightarrow^* , \sqsubseteq^* and \rightarrow^*) of Definition 2 are then computed considering the mutual dependency relations between them. Since we are only interested in reference types, we do not carry out analysis on basic types such as `boolean`, `int` and `double`. We also do not consider more advanced Java features such as functional interfaces and lambda expressions, as well as usages of Java Native Interface (JNI), nor method calls via Java reflective API. We have not tried to apply the approach to Java libraries, all invocation of methods from JDK are treated as end points, thus all possible call back edges will be missed in the analysis. Array accesses are treated conservatively—all type information that flows to one member of a reference array flows to all members of that array, so that only one node is generated for each array.

Since call graph information may be saved and be used for subsequent analyses, we propose the follow two ways to reduce storage of computed result. If two variables are similar regarding type information in a graph representation, they may be merged and then referred to by the merged node.

1. If $x \sqsubseteq^* y$ and $y \sqsubseteq^* x$, we say x and y form an *alias pair*, written $x \sim y$. Intuitively, VAR/\sim is a partition of VAR such that each $c \in \text{VAR}/\sim$ is a strongly connected component (SCC) in the variable graph edged by relation \sqsubseteq , which can be quickly collected by using Tarjan’s algorithm [?].
2. A more aggressive compression can be achieved in a way similar to bisimulation minimization of finite state systems [?,?]. Define $\approx \subset \text{VAR} \times \text{VAR}$ such that $x \approx y$ is symmetric and if
 - for all class c , $c \dashrightarrow^* x$ iff $c \dashrightarrow^* y$, and
 - for all $x \xrightarrow{f}^* x'$ there exists $y \xrightarrow{f}^* y'$ and $x' \approx y'$.

It is straightforward to see that \approx is a more aggressive merging scheme.

Lemma 2. *For all $x, y \in \text{VAR}$, $x \sim y$ implies $x \approx y$.*

We implement the second storage minimization scheme by using Kanellakis and Smolka’s algorithm [?] which computes the largest bisimulation relation for a given finite state labelled transition system. In our interpretation, the variables are treated as states and the field access relation is treated as a transition relation. The algorithm then merges equivalent variables into a single group. As a storage optimization process, this implementation has been tested and evaluated in the following section.

4 Experiment and Evaluation

We evaluate our approach by measuring its performance in terms of generated call sites and runtime on 13 benchmark programs. Among the benchmark programs, *compress*, *crypto* are from the SPECjvm2008 suite [?], and 11 programs are from the DaCapo suite [?]. We randomly selected these test cases from the two benchmark suites, in order to test code bases that are representative from a variety of sizes. All of our experiments were conducted on a Huawei Laptop equipped with Intel i5-8250U processor at 1.60 GHz and 8 GB memory, running Ubuntu 16.04LTS with OpenJDK 1.8.0.

We compare our approach against the default implementation of Class Hierarchy Analysis (CHA) and context-insensitive points-to analysis [?] that are implemented by the Soot team. We use Soot as our basic framework to extract the SSA based representation of the benchmark code. We also generate automata representation for the resulting relations which can be visualized in a subsequent user-friendly manual inspection. The choice of the context-insensitive points-to analysis is due to our approach also being context-insensitive, thus the results will be comparable. In the following tables we use CHA, PTA and TFA to refer to the results related to Class Hierarchy Analysis, Points-to Analysis and Type Flow Analysis, respectively.

During the evaluation the following three research questions are addressed.

| bench | $T_{\text{CHA}}(s)$ | $T_{\text{PTA}}(s)$ | $T_{\text{TFA}}(s)$ | $R_{\rightarrow\rightarrow}$ | R_{\sqsubseteq} | R_{\rightarrow} |
|--------------------|---------------------|---------------------|---------------------|------------------------------|-------------------|-------------------|
| compress | 0.02 | 0.11 | 0.01 | 87 | 202 | 24 |
| crypto | 0.01 | 0.13 | 0.04 | 94 | 226 | 18 |
| bootstrap | 23.97 | 34.59 | 0.05 | 191 | 453 | 17 |
| commons-codec | 0.008 | 0.12 | 0.14 | 306 | 3324 | 49 |
| junit | 24.56 | 34.01 | 0.19 | 1075 | 5772 | 241 |
| commons-httpclient | 0.008 | 0.11 | 0.38 | 2423 | 8511 | 521 |
| serializer | 22.95 | 32.72 | 1.64 | 3006 | 17726 | 331 |
| xerces | 22.49 | 31.18 | 2.71 | 12590 | 72503 | 2779 |
| eclipse | 22.80 | 39.74 | 1.64 | 7933 | 37435 | 1620 |
| derby | 22.71 | 49.05 | 18.13 | 20698 | 191854 | 5386 |
| xalan | 79.57 | 163.68 | 35.69 | 32971 | 162249 | 3696 |
| antlr | 43.08 | 92.18 | 3.84 | 16117 | 76741 | 3879 |
| batik | 48.85 | 94.20 | 7.19 | 29409 | 122534 | 6039 |

Fig. 7. Runtime cost with different analysis

RQ 1: How efficient is our approach compared with the traditional Class Hierarchy Analysis and points-to analysis?

RQ 2: How accurate is the result of our approach when comparing with the other analyses?

RQ 3: Does our optimization algorithm achieve significantly reduced space consumption?

4.1 RQ1: Efficiency

To answer the first research question, we executed each benchmark program 10 times with the CHA, PTA and TFA algorithms. We calculated the average time consumption as displayed at the left three columns of the Table in Figure 7. The sizes of each relation that our approach generated (i.e., the type flow relation ‘ $\rightarrow\rightarrow$ ’, variable partial order ‘ \sqsubseteq ’ and the field access ‘ \rightarrow ’) are counted, which provides a feeling of size for the problem we are treating. One may observe that when the problem size increases, the execution of our algorithm also increases in a way similar to CHA, though in general the runtime of CHA is supposed to grow linearly in the size of a program. The reason that TFA sometimes outperforms CHA may be partially due to the intrinsic complexity of the class and interface hierarchy that a program adopts. TFA is in general more efficient than PTA. The runtime cost in TFA basically depends on the size of generated relations, as well as the relational complexity as most of the time is consumed to calculate a fixpoint. For PTA it requires extra time for maintaining a heap abstraction.

Taking a closer look at the benchmark *bootstrap*, CHA, PTA analyze the benchmark using about 23.74 and 34.13 seconds, respectively. As TFA only generated 661 relations, the analysis only takes 0.03 second. There are 3 benchmarks, *xalan*, *antlr*, *batik*, that cannot be analyzed by PTA, marked as “-” in the table in

| bench | CS_{base} | CS_{CHA} | CS_{PTA} | CS_{TFA} |
|--------------------|--------------------------|-------------------------|-------------------------|-------------------------|
| compress | 153 | 160 | 18 | 73 |
| crypto | 302 | 307 | 62 | 121 |
| bootstrap | 657 | 801 | 891 | 328 |
| commons-codec | 1162 | 1372 | 270 | 554 |
| junit | 3196 | 17532 | 11176 | 1197 |
| commons-httpclient | 6817 | 17118 | 567 | 2928 |
| serializer | 4782 | 9533 | 1248 | 1756 |
| xerces | 24579 | 56252 | 10631 | 8120 |
| eclipse | 23607 | 95073 | 70016 | 9409 |
| derby | 69537 | 180428 | 85160 | 16432 |
| xalan | 57430 | 155866 | 164805 | 18682 |
| antlr | 62007 | 147014 | 105861 | 17187 |
| batik | 56877 | 235085 | 145126 | 20990 |

Fig. 8. Call sites generated by different analyses

Figure 7. The points-to analysis has caused an exception “OutOfMemoryError with JVM GC overhead limit exceeded”, as the JVM garbage collector is taking an excessive amount of time (by default 98% of all CPU time of the process) and recovers very little memory in each run (by default 2% of the heap).

4.2 RQ2: Accuracy

We answer the second question by considering the number of generated call sites as an indication of accuracy. In type flow analysis, a method call $a.m()$ is resolved to $c.m()$ if class c is included in a ’s reaching type set and method $m()$ is defined for c . In general, a more accurate analysis often generates a smaller set of types for each calling variable, resulting fewer call edges in total in the call graph. The table included in Figure 8 displays the number of call sites generated by different analyses. We also include the origin call site counting, i.e., the number of call sites syntactically written in the source code, as the baseline at the **CS_{base}** column. In comparison to CHA, our approach has reduced a significant amount of call site edges. Comparing to other two analyses, the number of call edges resolved by TFA are often larger than PTA and smaller than CHA on the same benchmark. The difference may be caused by our over approximation on analyzing array references, as well as the existence of unsolved call edges from e.g. JNI calls or reflective calls.

4.3 RQ3: Optimization

We apply bisimulation minimization to merge nodes that have the same types as well as accessible types through fields. Thus we can reduce the space consumption when storing the result for subsequent analysis processes. Regarding the third research question, we calculated the number of relation nodes before and after optimization. Besides, time consumption is another factor that we consider.

| bench | Node _{origin} | Node _{opt} | Reduce | Time(s) |
|--------------------|------------------------|---------------------|--------|---------|
| compress | 205 | 130 | 36.59% | 0.008 |
| crypto | 312 | 158 | 49.36% | 0.010 |
| bootstrap | 514 | 279 | 45.72% | 0.019 |
| commons-codec | 1742 | 886 | 49.14% | 0.202 |
| junit | 5859 | 3243 | 44.65% | 2.269 |
| commons-httpclient | 9708 | 5164 | 46.81% | 4.651 |
| serializer | 9600 | 6668 | 30.54% | 3.198 |
| xerces | 41634 | 25351 | 39.11% | 912.80 |
| eclipse | 34631 | 19170 | 44.64% | 78.26 |
| derby | 112265 | - | - | - |
| xalan | 103697 | - | - | - |
| antlr | 56589 | - | - | - |
| batik | 89336 | - | - | - |

Fig. 9. Optimization result

The results are shown in the table in Figure 9. We evaluated our optimization algorithm on all benchmarks 10 times and successfully executed 7 out of the 13 benchmarks. For those succeed, we can reduce the space by about 45% on average, with a reasonable time consumption. The other 6 benchmarks cannot be due to the exception “OutOfMemoryError with JVM GC overhead limit exceeded” being thrown, for which we leave the symbol “-” accordingly in the table.

5 Related Work

There are not many works focusing on general purpose call graph construction algorithms, and we give a brief review of these works first.

As stated in the introduction, Class Hierarchy Analysis (CHA) [?,?], Rapid Type Analysis (RTA) [?] and Variable Type Analysis (VTA) [?] are efficient algorithms that conservatively resolves call sites without any help from points-to analysis. Grove et al. [?] introduced an approach to model context-sensitive and context-insensitive call graph construction. They define call graph in terms of three algorithm-specific parameter partial orders, and provide a method called Monotonic Refinement, potentially adding to the class sets of local variables and adding new contours to call sites, load sites, and store sites. Tip and Palsberg [?] Proposed four propagation-based call graph construction algorithms, CTA, MTA, FTA and XTA. CTA uses distinct sets for classes, MTA uses distinct sets for classes and fields, FTA uses distinct sets for classes and methods, and XTA uses distinct sets for classes, fields, and methods. The constructed call graphs tend to contain slightly fewer method definitions when compared to RTA. It has been shown that associating a distinct set of types with each method in a class has a significantly greater impact on precision than using a distinct set for each field in a class. Reif et al. [?] study the construction of call graphs for Java libraries that abstract over all potential library usages, in a so-called *open world* approach. They invented two concrete call graph algorithms for libraries

based on adaptations of the CHA algorithm, to be used for software quality and security issues. In general they are interested in analyzing library without knowing client application, which is complementary to our work that has focus on client program while treating library calls as end nodes.

Call graphs may serve as a basis for points-to analysis, but often a points-to analysis implicitly computes a call graph on-the-fly, as the context insensitive points-to algorithm implemented in Soot using SPARK [?]. Most context-sensitive points-to analysis algorithms [?,?,?] progress call edges at the same time with value flow, to our knowledge. The main distinction of our approach from these points-to analysis is the usage of an abstract heap, as we are only interested in the actual reaching types of the receiver of a call. However, unlike CHA and VTA, our methodology can be extended to context-sensitive settings.

6 Conclusion

In this paper we have proposed Type Flow Analysis (TFA), an algorithm that constructs call graph edges for Object-Oriented programming languages. Different from points-to based analysis, we do not require a heap abstraction, so the computation is purely relational. We have proved that in the context-insensitive setting, our result is equivalent to that would be produced by a subset-based points-to analysis, regarding the core Object-Oriented language features. We have implemented the algorithm in the Soot compiler framework, and have conducted preliminary evaluation by comparing our results with those produced by the built-in CHA and Points-to analysis algorithms in Soot on a selection of 13 benchmark programs from SPECjvm2008 and DaCapo benchmark suites, and achieved promising results. In the future we plan to develop context-sensitive analysis algorithms based on TFA.