

# Data Flow Analysis in the Presence of Correlated Calls

Marianna Rapoport<sup>1</sup>, Ondřej Lhoták<sup>1</sup>, and Frank Tip<sup>2</sup>

<sup>1</sup> University of Waterloo  
{mrapoport, olhotak}@uwaterloo.ca  
<sup>2</sup> Samsung Research America  
ftip@samsung.com

Frank wanted to discuss the title (M)

**Abstract.** We present a technique to improve the precision of data-flow analyses on object-oriented programs in the presence of *correlated calls*. Two method calls are correlated if they are polymorphic and are invoked on the same object. Correlated calls are problematic because they can make existing data-flow analyses consider certain infeasible data-flow paths as valid. This leads to loss in precision of the analysis solution. We show how infeasible paths can be eliminated for *Inter-procedural Finite Distributive Subset* (IFDS) problems, a large class of data-flow analysis problems. We show how the precision of IFDS problems can be improved in the presence of correlated calls, by using the *Inter-procedural Distributive Environment* (IDE) algorithm to eliminate infeasible paths. Using IDE, we eliminate the infeasible paths and obtain a more precise result for the original IFDS problem. Our analysis is implemented in Scala, using the WALA framework for static program analysis on Java bytecode.

## 1 Introduction

Static program analysis aims to discover properties of computer programs without running them. Static analysis has applications in compiler optimization, development of programming tools, and computer security, among others. As an example, we might want to analyze a program to know which variables are constants. We could then write a compiler optimization that ensures that the values of those variables are computed only once. Alternatively, we could use the information about constant variables in an integrated development environment; for instance, to notify the user when an if-expression executes only one of its branches because its test condition has a constant value.

There are demonstrable limits on what information we can obtain about a program without running it. Rice's theorem states that verifying any non-trivial property of a program is an undecidable problem [18]. However, it is sometimes possible to design an algorithm that *over-* or *underapproximates* the solution that we are seeking.

This section should be shrunk a lot (not completely deleted) (O)

*Data-flow analysis* is an area of program analysis whose goal is to compute approximations of certain information (for example, which variables must be constants) for each program point.

Other examples of data-flow analyses are *reaching definitions* (finding out up to which instruction a given assignment of a variable must be valid) and *available expressions* (retrieving the expressions in the program that do not need to be recomputed at a given program point).

Another example of a data-flow analysis is *taint analysis* [24]. Taint analysis discovers if “secret” values, like passwords or other confidential user information, can leak to an external observer. Methods that generate secret values, e.g. those that read user input, are called *sources*. Methods that can leak information, e.g. those that write data to a file or send data through a network, are called *sinks*. The goal of taint analysis is to find out whether data can propagate from sources to sinks.

An important property of a data-flow analysis is *precision*. Precision reflects how closely a data-flow-analysis result over- or underapproximates the information we are interested in. In the case of taint analysis, let  $T$  be the number of sinks that the analysis considers to leak secret information, and  $R$  the real number of potential information leaks. The smaller the difference between  $T$  and  $R$ , the greater the precision of the taint analysis.

Data-flow analyses operate on *control-flow graphs* that model the order in which the instructions of a program are executed. A data-flow-analysis problem defines *flow functions* that represent how data is propagated along the edges of the control-flow graph. The *confluence operator* specifies how the data that has been computed along different paths should be merged when the paths join.

Since a control-flow graph is an overapproximation of the possible flows of control in concrete executions of a program, the graph may contain *infeasible* paths that cannot occur at runtime.

One way to improve the precision of a data-flow analysis is to detect and eliminate infeasible paths.

Our goal is to improve the precision of solutions to problems that can be solved by the *Inter-procedural Finite Distributive Subset* (IFDS) algorithm [17]. The IFDS algorithm is a general data-flow algorithm that can compute solutions to various data-flow problems, like reaching definitions, available expressions, and taint analysis.

We improve the precision of IFDS problem solutions by eliminating infeasible paths that occur in object-oriented programs in the presence of *correlated method calls* — polymorphic calls that are invoked on the same object [?].

## 2 Motivating Example

Consider a call site  $r.m()$  in an object-oriented programming language, where the variable  $r$  is the *receiver* variable of the call site and  $m$  is the name of the

invoked method<sup>3</sup>. In the rest of the paper, we use the general term *receiver* to mean a receiver variable. At runtime, the actual method that will be invoked by the call site depends on the runtime type of the object referenced by  $r$ . If the call site  $r.m()$  can be associated with more than one method at compile time, we will say that the call site is *polymorphic*.

For example, in Listing 1, it is not possible to infer statically whether the runtime type of the variable `a` in the `main` method is `A` or `B`. The call `a.foo()` can be dispatched to either `A.foo` or `B.foo`, and `a.bar(v)` can be dispatched to either `A.bar` or `B.bar`. A concrete execution path for the `main` method might therefore go through `A.foo` and `A.bar`, or through `B.foo` and `B.bar`. However, there cannot be an execution path through `A.foo` and `B.bar` or through `B.foo` and `A.bar`.

```

class A {
    String foo {
        return secret();
    }

    void bar(String s) {}
}

class B extends A {
    String foo {
        return "not_secret";
    }

    void bar(String s) {
        System.out.println(s);
    }
}

class Main {
    public static void main(String[] args) {
        A a = args == null ? new A() : new B();
        String v = a.foo();
        a.bar(v);
    }
}

```

**Fig. 1:** Example program containing correlated calls

---

<sup>3</sup> We assume an internal representation of the program in which for each call site  $e_r.m()$ , the expression  $e_r$  has been evaluated to the variable  $r$ .

We call the invocations to methods `foo` and `bar` *correlated*. More generally, correlated calls occur when more than one polymorphic call is invoked on the same receiver variable.

Suppose we wanted to perform a taint analysis on the program in Listing 1. Most dataflow-analysis algorithms, including IFDS, would conservatively assume that the call `a.bar` could be dispatched to both `A.bar` and `B.bar`, independently of what `a.foo` had been dispatched to in the previous line.

As a result, such an analysis would consider a path through `A.foo` and `B.bar` feasible. This means that the variable `v` would be considered secret. We would conclude that a secret value is passed to `B.bar` and printed to the user. In other words, we would consider the program to leak secret information, which it does not do in any concrete execution.

Our technique for improving the precision of an IFDS result is based on transforming the original IFDS problem into a more expressive *Inter-procedural Distributive Environment* (IDE) problem. IDE problems can be solved with the IDE algorithm which is a generalization of IFDS [20]. The IDE algorithm can, for instance, solve certain versions of the constant propagation problem that IFDS cannot.

To improve the precision of IFDS results, given an IFDS problem  $P$ , we convert it into an IDE problem  $Q$  that accounts for correlated calls. We then use the IDE algorithm to obtain a solution to  $Q$ . Finally, we convert the IDE result into a IFDS result. In the presence of correlated calls, the obtained IFDS result can be more precise than the solution that the IFDS algorithm would compute for  $P$ .

The goal of the correlated-calls analysis presented in this work is to modify the output of an IFDS analysis to account for correlated calls. Specifically, the correlated-calls analysis improves the precision of IFDS problem results by eliminating infeasible execution paths caused by correlated calls. This is done by converting the input-IFDS problem to an IDE problem that detects infeasible paths, and converting the IDE result back to a more precise IFDS result.

The contributions of this work are:

- A transformation from IFDS to IDE problems that considers correlated calls.
- An implementation in Scala of the correlated-calls transformation and the IDE algorithm which is based on the WALA framework for static analysis on Java bytecode [6].

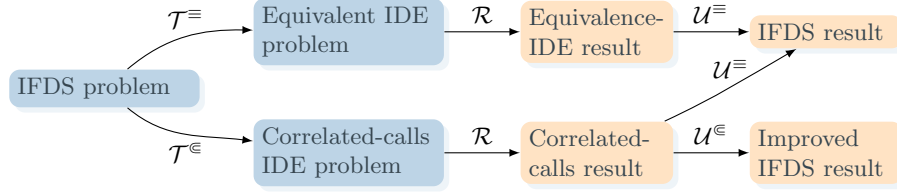
Reference the figure here. (M)

We prove that the solution to an IDE problem that considers correlated calls is more precise than the solution to the original IFDS problem. We also show that the correlated-calls analysis is sound, i.e. that it never considers concrete execution paths as infeasible.

As discussed on the phone, we want to talk more about the results here. (M)

Finally, we evaluate the effectiveness of the correlated-calls analysis using an implementation of taint analysis as the source IFDS problem.

The remainder of this paper is organized as follows. In the next section, we describe the IFDS and IDE analyses in detail. In Section 4 we present the



**Fig. 2:** Transformations between IFDS and IDE problems and their results

correlated-calls analysis as a transformation of IFDS problems into a special kind of IDE problem. Section 5 describes an efficient representation of the data structures that are required to define a correlated-calls IDE transformation. In Section 6 we address some implementation aspects of the correlated-calls analysis and present an evaluation of its results. Section 8 contains concluding remarks. The proofs for the lemmas in the paper are provided in the Appendix.

Reference appendix like this? (M)

### 3 Background

The purpose of the correlated-calls analysis is to solve IFDS problems more precisely than using the standard IFDS algorithm by ruling out some infeasible paths. The correlated-calls analysis works by transforming an IFDS problem to an IDE problem, solving the IDE problem, and transforming the IDE result to a solution to the original IFDS problem. This section describes the general ideas underlying IFDS and IDE.

#### 3.1 Terminology and Notation

We will start by introducing several concepts used by the IFDS and IDE analyses.

A *control-flow graph* is a directed graph in which nodes correspond to instructions and edges represent transfer of control between the instructions during an execution of the program. A control-flow graph has a unique start node,  $\text{start}_{\text{main}}$ , which is the node corresponding to the program entrypoint.

An *intra-procedural* path is a path in a control-flow graph whose nodes are in the same procedure. By contrast, an *inter-procedural path* is one that contains nodes from different procedures.

A *control-flow supergraph* is a control-flow graph in which each procedure  $p$  is augmented with an additional *start node*  $\text{start}_p$  and *end node*  $\text{end}_p$ , and for each call  $c_q$  to a procedure  $q$ , there is a *call node*  $\text{call}_{c_q}$  and subsequent *return node*  $\text{return}_{c_q}$ .

A control-flow supergraph allows us to model the control flow in inter-procedural paths. The flow from the caller to the callee is represented using an edge

$$(\text{call}_{c_q}, \text{start}_q).$$

The control flow from the callee back to the caller goes through an edge

$$(\text{end}_q, \text{return}_{c_q}).$$

A *flow-sensitive* data-flow analysis is one that takes the order of program instructions into account.

Let each call node in a program be labeled with a distinct opening parenthesis and the corresponding return node with the matching closing parenthesis. For a given path  $p$ , let  $s$  be the string that is obtained by concatenating the labels of the nodes in  $p$ . Then  $p$  is *valid* if  $s$  belongs to the language of substrings of balanced parentheses. The set of all inter-procedurally valid paths from the start node to a node  $n$  is denoted as  $\text{VP}(n)$ . The set  $\text{VP}(n)$  is a conservative approximation of all concrete execution paths from the start node to  $n$ .

A *context-sensitive* data-flow analysis is an analysis that considers only inter-procedurally valid paths.

*Example 1.* In the supergraph in Figure ??, let us assign  $\{, \}$  parentheses to  $\text{call}_{A.f}$  and  $\text{return}_{A.f}$ , and  $\langle, \rangle$  parentheses to  $\text{call}_f$  and  $\text{return}_f$ . Then the string corresponding to the path

$$p_1 = [\text{call}_{A.f}, \text{start}_f, \text{if } (s == \text{null}), \text{return } s, \text{end}_f, \text{return}_{A.f}]$$

is  $\{\}$ , which indicates that  $p_1$  is valid. Every prefix of  $p_1$  is also a valid path.

However, the graph also contains an inter-procedurally invalid path

$$p_2 = [\text{call}_f, \text{start}_f, \text{if } (s == \text{null}), \text{return } s, \text{end}_f, \text{return}_{A.f}]$$

with corresponding string  $\langle \rangle$ .

A *lattice* is a partially ordered set in which each subset has a least upper bound and a greatest lower bound.

A *meet semilattice*  $L = (S, \sqcap)$  is defined by a set  $S$  and a meet operation  $\sqcap$  that is associative, commutative, and idempotent. The meet operation induces a partial order  $(S, \sqsubseteq)$  where every subset contains a greatest lower bound: For all  $x, y \in S$ ,  $x \sqsubseteq y$  if  $x \sqcap y = x$ . The greatest lower bound, or top element, of the semilattice is denoted as  $\top$ . If  $k$  is the length of the longest chains of elements in the semilattice, then the *height* of the semilattice is  $k - 1$ .

Finally, we introduce the notion of *distributivity*. Given a set  $D$ , a function  $f : 2^D \rightarrow 2^D$  is distributive if  $\forall x_1, x_2 \in 2^D$ ,

$$f(x_1 \cup x_2) = f(x_1) \cup f(x_2). \quad (1)$$

In this work, we will denote a map from a set of keys  $K$  to values from set  $V$  as

$$\{(k, v) \mid k \in K, v \in V\}. \quad (2)$$

For an arbitrary map  $m$ ,  $m(x)$  is the value to which  $x$  is mapped in  $m$ . We denote by  $m[x \rightarrow y]$  a map identical to  $m$ , except that the element  $x$  is mapped to  $y$ . To avoid excessive parentheses, we write  $(m[x_1 \rightarrow y_1])[x_2 \rightarrow y_2]$  as  $m[x_1 \rightarrow y_1][x_2 \rightarrow y_2]$ .

We will denote the identity function  $\lambda x . x$  by  $\text{id}$ . We will use a typed version of this function in various contexts, where the type of  $x$  will vary with the context.

The figure that we deleted was constructed so that there would be a realizable and a non-realizable path. Now that we deleted that figure, let's delete this example as well. (M)

### 3.2 IFDS

The IFDS framework is a precise and efficient algorithm for data-flow analysis. IFDS was developed in 1995 by T. Reps, S. Horwitz, and M. Sagiv at the University of Wisconsin and has been used to solve a variety of data-flow analysis problems [4,13,10,24].

The IFDS algorithm is applicable to problems which can be expressed with data-flow functions that satisfy certain restrictions. *Inter-procedural* flow functions specify how data flows from the invocation of a procedure to its start, and from the procedure’s end back to its call site. *Distributive* flow functions are those that distribute over the confluence operator. In the context of IFDS, the confluence operator is called *meet*, and it can be either union or intersection. The data-flow facts on which the analysis operates must be a *finite* set  $D$ . Each flow function operates on a *subset* of  $D$  (for example, the set of variables in the program) which makes the domain of the flow functions the power set of  $D$ .

Given a data-flow problem that satisfies the restrictions of IFDS, the algorithm provides a *context-sensitive* solution in polynomial time. Compared to IFDS, most data-flow analyses are either general but do not run in polynomial time [8,21] or handle a very specific set of problems [9].

The four references take away space in the References section. Remove some of them? (M)

I’m not saying any more that IFDS is a version of Sharir&Paueli’s functional approach. (M)

I also removed the comparison of IFDS with other algorithms. I’m not pointing out the benefits of IFDS any more, so we might fail to show the importance of the algorithm. (M)

Provide reference explaining context sensitivity? (M)

**Data-Flow Problems Suitable for IFDS** Informally, an IFDS analysis can only solve decision problems. An IFDS analysis answers questions of the following kind: “is property  $X$  true at program point  $Y$ ?”. For example, a taint-analysis problem asks, for each variable  $v$  in the program, “is  $v$  secret at a given program point?”.

Formally, a data-flow analysis problem is suitable for an IFDS analysis if it can be encoded as an IFDS problem

$$(G^*, D, F, M_F, \sqcap), \quad (3)$$

where  $G^* = (N^*, E^*)$  is the supergraph of the input program with nodes  $N^*$  and edges  $E^*$ ,  $D$  is a finite set of *data-flow facts*,  $F$  is a set of distributive dataflow functions of type  $2^D \rightarrow 2^D$ ,  $M_F : E^* \rightarrow F$  is a function that maps supergraph edges to dataflow functions, and  $M_F$  is extended to paths by composition<sup>4</sup>. The *meet operator*  $\sqcap$  is either union or intersection.

Without loss of generality, we will take *meet* to denote union. It can be shown that any problem where *meet* is defined as intersection can be reformulated into an equivalent one where *meet* is defined as union [17].

**Overview of the IFDS Algorithm** Given an IFDS problem, for each node  $n \in N^*$  the IFDS algorithm computes the *meet-over-all-valid-paths* solution

$$\text{MVP}_F(n) = \bigsqcap_{q \in \text{VP}(n)} M_F(q)(\emptyset). \quad (4)$$

<sup>4</sup> Let  $A$  be a set and  $f : E^* \rightarrow (A \rightarrow A)$  a function from supergraph edges to functions on  $A$ . We say that  $f$  is extended to paths by composition to denote that for a path  $q$  consisting of the edges  $e_1, \dots, e_k$ ,  $f(q) = f(e_k) \circ \dots \circ f(e_1) \circ \text{id}$ .

To compute the meet-over-all-valid-paths solution, each node in the control-flow supergraph is paired with a *fact*  $d \in D \cup \{\mathbf{0}\}$ ,  $\mathbf{0} \notin D$ , yielding the nodes  $N^\#$  of the *exploded supergraph*  $G^\# = (N^\#, E^\#)$ . Roughly, for each node in the program, a fact denotes a binary property whose value (true or false) we want to find out. The start node of the exploded supergraph is the node  $(\text{start}_{\text{main}}, \mathbf{0})$ .

The flow functions  $F$  define the edges of the exploded supergraph. Using the flow functions, the IFDS algorithm computes the inter-procedurally *realizable* paths from the start to the rest of the exploded graph's nodes. A *realizable* path is a valid path in the exploded supergraph that starts with the entry node  $\text{start}_{\text{main}}$ .

If there is a realizable path from the node  $(\text{start}_{\text{main}}, \mathbf{0})$  to a given node  $(n, d)$ ,  $d \neq \mathbf{0}$ , then the fact  $d$  is considered to hold at node  $n$ . A path to a node  $(n, \mathbf{0})$  means that in the control-flow supergraph, there is a path from  $\text{start}_{\text{main}}$  to  $n$ .

In this way, the IFDS algorithm reduces the input data-flow problem to a graph-reachability problem.

*Example 2.* In a taint analysis,  $D$  is the set of variables in the program. If a fact  $d \in D$  is reachable at a given node, then the variable is considered secret at that node. Otherwise, it is considered not secret. The question “is  $d$  secret at node  $n$ ?” becomes “is there a realizable path from  $(\text{start}_{\text{main}}, \mathbf{0})$  to  $(n, d)$ ?”.

*Example 3.* The exploded supergraph for Listing 1 is shown in Figure 3. We can see that there is a realizable path from the start node of the exploded graph to the variable  $\mathbf{s}$  at the node `print(s)` in the `B.bar` method. This means that at that node,  $\mathbf{s}$  is considered secret.

The flow functions  $F \subseteq 2^D \rightarrow 2^D$  allow us to establish the edges in the exploded supergraph.

Given a control-flow-graph edge  $e = (n_1, n_2) \in E^*$  and a distributive dataflow function  $f = M(e)$ , the *representation relation*  $R_f : (D \cup \{\mathbf{0}\}) \times (D \cup \{\mathbf{0}\})$  of  $f$  is defined as

$$R_f = \{(\mathbf{0}, \mathbf{0})\} \cup \{(\mathbf{0}, d_j) \mid d_j \in f(\emptyset)\} \cup \{(d_i, d_j) \mid d_j \in f(\{d_i\}), d_j \notin f(\emptyset)\}.$$

Each pair  $(d_i, d_j) \in R_f$  corresponds to an edge  $((n_1, d_i), (n_2, d_j))$  in the exploded supergraph.

Note that  $R_f$  constructs pairs of dataflow facts so that

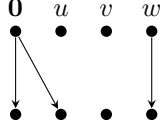
- there is always an edge  $(\mathbf{0}, \mathbf{0})$  corresponding to the control-flow-graph edge;
- if there is an edge  $(\mathbf{0}, d_j)$ , then there is no other edge leading to  $d_j$ ; in particular, there is never an edge  $(d_i, \mathbf{0})$  where  $d_i \neq \mathbf{0}$ .

*Example 4.* The representation relation  $R_f$  for a set of data-flow facts  $D = \{u, v, w\}$  and dataflow function  $f = \lambda S. S \setminus \{v\} \cup \{u\}$  looks as follows:

$$R_f = \{(\mathbf{0}, \mathbf{0}), (\mathbf{0}, u), (w, w)\}.$$

The corresponding exploded-graph edges are shown below.





The representation relation lets us decompose a flow function into functions that operate on each fact individually. This is possible due to distributivity: we can apply the flow function on each single fact and take the union of the results, rather than applying the function to the union of the facts.

The representation relation allows us to compactly represent the composition and meet operations which are required for the IFDS algorithm.

For two representation relations  $R_{f_1}$ ,  $R_{f_2}$ , the composition and meet operations are defined as follows:

$$R_{f_1} \circ R_{f_2} = \{(d_1, d_3) \mid \exists d_2 : (d_1, d_2) \in R_{f_1}, (d_2, d_3) \in R_{f_2}\}. \quad (5)$$

and

$$R_{f_1} \sqcap R_{f_2} = R_{f_1} \cup R_{f_2}. \quad (6)$$

The representation relation distributes over composition and meet:

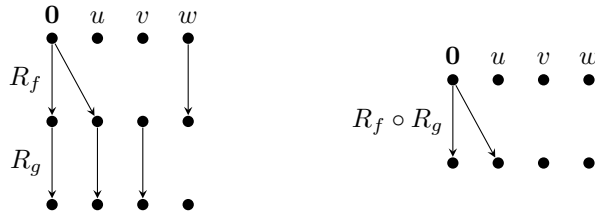
$$R_{f_1} \circ R_{f_2} = R_{f_1 \circ f_2}, \quad \text{and} \quad R_{f_1} \sqcap R_{f_2} = R_{f_1 \sqcap f_2}. \quad (7)$$

On the exploded graph, the composition of two functions is represented by the paths that are formed when the exploded-graph edges are combined.

*Example 5.* If  $g = \lambda S. S \setminus \{w\}$  and  $f$  is defined as in example 4, then

$$R_f \circ R_g = \{(\mathbf{0}, \mathbf{0}), (\mathbf{0}, u)\},$$

as illustrated by the corresponding exploded graph edges:



We presented an overview of the IFDS analysis. IFDS problems are transformed into IDE problems by the correlated-calls analysis. The IDE framework is described in the next section.

We never use the interpretation function, so I deleted it, too (M)

### 3.3 IDE

There exists an entire class of data-flow problems that cannot be formulated as IFDS problems. Informally, the problems cannot be formulated as decision problems. For instance, a constant-propagation problem asks, for each variable  $v$  in the program, “if  $v$  is a constant at a given program point, what is  $v$ ’s value?”. The questions asked by constant propagation are of the form “if property  $X$  ( *$v$  being a constant*) is true at program point  $Y$ , what is the value of some property  $Z$  (*the value of the constant*) corresponding to  $X$ ?”. It turns out that problems with such questions can often be solved by the IDE algorithm.

The IDE framework is an expressive extension to IFDS that was created by the same authors in 1996. The problems that IDE can solve include, but are not limited to, IFDS problems [17]. Just as the IFDS algorithm, the IDE algorithm is suitable for data-flow analyses that can be encoded with inter-procedural, distributive flow functions. However, in IDE, the domain of the flow functions is not restricted to sets  $D$  of data-flow facts. The IDE domain of a flow function consists of *environments* that map data-flow facts from the set  $D$  to lattice elements.

As an example, in a constant propagation problem, an IDE environment would map each variable to the (possibly) constant value that it is bound to. To illustrate the distinction between IFDS and IDE we could say that IFDS can find out which variables in a program are constants, whereas IDE can additionally retrieve the values of the constant variables. Instead of just telling us whether a fact holds or not, the IDE analysis can provide us with additional information about facts.

Just as in the IFDS analysis, the IDE algorithm reduces a data-flow problem to a graph-reachability problem. Additionally, for each program point, the algorithm computes an *environment*  $\text{Env}(D, L)$ , where data-flow facts are mapped to values of a lattice  $L$ .

For example, using the IDE analysis, we can encode a restricted version of a constant-propagation analysis<sup>5</sup>. The data-flow facts correspond to program variables, and the lattice incorporates all possible values for constants. If a fact  $d$  in the exploded supergraph is reachable at node  $n$ , and  $\text{Env}(d) \notin \{\perp, \top\}$ , it means that the variable associated with  $d$  is a constant. Furthermore, the value of the constant can be inferred from the environment for the corresponding node and is equal to  $\text{Env}(d)$ .

Formally, an IDE problem is defined as a four-tuple

$$(G^*, D, L, M_{\text{Env}}), \quad (8)$$

where  $G^*$  is a control-flow supergraph,  $D$  is a set of data-flow facts, and  $L$  is a meet semilattice with finite height. Finally,  $M_{\text{Env}} : E^* \rightarrow (\text{Env}(D, L) \rightarrow \text{Env}(D, L))$  is a function from the edges of the control-flow supergraph to distributive *environment transformers*.  $M_{\text{Env}}$  is extended to paths by composition.

<sup>5</sup> In the general case, constant propagation cannot be encoded with distributive flow functions and is therefore not suitable for an IDE analysis [12].

Given an IDE problem, for each node  $n \in N^*$  and fact  $d \in D$ , the IDE algorithm computes the meet-over-all-valid-paths solution

$$\text{MVP}_{\text{Env}}(n, d) = \bigcap_{q \in \text{VP}(n)} M_{\text{Env}}(q)(\Omega)(d), \quad (9)$$

where  $M_{\text{Env}}$  is extended to paths by composition and

$$\Omega = \lambda d. \top \quad (10)$$

is the top element in the environment lattice  $\text{Env}(D, L)$ .

**Environment Transformers** For each node in the control-flow graph, the result of an IDE analysis computes an environment  $\text{Env}(D, L)$ , which is a map from data-flow facts to lattice elements.

Instead of flow functions that show how to propagate facts, the IDE framework uses distributive environment transformers to propagate environments. For each edge  $(n_1, n_2)$  in the control-flow supergraph, an environment transformer indicates how the environment at node  $n_1$  is modified at node  $n_2$ .

From Section 3.2 we know that flow functions can be represented with exploded-graph edges. To represent environment transformers, we will construct *labeled* exploded-graph edges, where each edge is associated with a distributive *micro function*<sup>6</sup>  $f : L \rightarrow L$ . A micro function shows how to change a lattice element for a given node and fact.

For instance, if an IDE problem is equivalent to an IFDS problem, the edges of the exploded supergraph are the same for both problems. In the IDE problem, the edges of the exploded supergraph are labeled with identity micro functions.

We extend the meet operator to work on micro functions by defining

$$(f_1 \sqcap f_2)(l) = f_1(l) \sqcap f_2(l) \quad (11)$$

for all  $l \in L$ .

In IDE problems, the auxiliary fact analogous to  $\mathbf{0}$  in IFDS is denoted as  $\Lambda$ .

In this way, each edge in the exploded graph is labeled with a micro function. The mapping from exploded-graph edges to the corresponding micro functions is stored in *edge functions*, denoted as  $\text{EdgeFn} : E^\# \rightarrow (L \rightarrow L)$ .

**Overview of the IDE Algorithm** Given a labeled exploded supergraph, the IDE algorithm computes the environments for all nodes in the control-flow graph.

The algorithm first computes the lattice elements  $l_{n,d}$  that correspond to each reachable node  $(n, d)$  in the exploded supergraph. The union of the exploded nodes  $(n, d)$  for a given control-flow node  $n$ , mapped to the corresponding lattice elements  $l_{n,d}$ , form the environment  $\text{Env}_n$  for that node:

$$\text{Env}_n = \{(d, l_{n,d}) \mid (n, d) \in N^\#\}. \quad (12)$$

The overall idea behind computing the lattice elements  $l_{n,d}$  is the following. For each inter-procedurally realizable path

$$p = [(\text{start}_{\text{main}}, A), (n_1, d_1), \dots, (n_k, d_k)]$$

that starts with the entrypoint of the exploded supergraph, we compute the micro function  $f_p$  that corresponds to  $p$ . The micro function consists of the composition of all individual micro functions with which the edges of  $p$  are labeled:

$$f_p = \text{EdgeFn}((n_{k-1}, d_{k-1}), (n_k, d_k)) \circ \dots \circ \text{EdgeFn}(\text{start}_{\text{main}}, A), (n_1, d_1)). \quad (13)$$

Let the lattice element that  $(n, d)$  is mapped to according to path  $p$  be denoted as  $l_{n,d}^p$ . As shown in Sagiv et al. [20], the lattice element can be obtained by applying  $f_p$  to the bottom element:

$$l_{n,d}^p = f_p(\perp). \quad (14)$$

Let  $Q$  be the set of paths that start at the entry point and end at the given node  $(n, d)$ . The lattice element  $l_{n,d}$  is the meet of the lattice elements corresponding to all the paths in  $Q$ :

$$l_{n,d} = \bigcap_{q \in Q} l_{n,d}^q.$$

This is a general outline of the IDE analysis. We use the IDE framework to improve the precision of IFDS problems in the presence of correlated calls. The next section describes how this is done.

we agreed to move this figure to the initial example section (O)

Move Figure 3 up to where it's mentioned first? If yes, the caption won't make sense. (M)

Can we combine the IDE Figure3 with the IFDS Figure 4? It should be the same figure, just that Figure3 has labels on the edges. We can initially tell the reader to ignore the labels until later. (O)

The following definitions are very general, and not related to the transformation. Can we move them somewhere earlier? (O)

*Edge Functions* Let  $\mathcal{F}$  be the set of methods in a program with a signature  $s_{\mathcal{F}}$ .

**Definition 1.** Let  $r.c()$  be a call site on a receiver  $r \in R$  with runtime type  $t \in T$ . Let  $s_{\mathcal{F}}$  be the method signature corresponding to the call  $c()$ . For  $s_{\mathcal{F}}$  and  $t$ , a lookup function returns the method implementation  $f \in \mathcal{F}$  to which the call  $r.c()$  is dispatched:

$$\text{lookup}(s_{\mathcal{F}}, t) = f. \quad (15)$$

**Definition 2.** For a method signature  $s_{\mathcal{F}}$  and a method implementation  $f \in \mathcal{F}$ , the static-type function  $\tau$  returns the set of types for which the lookup function yields  $f$ :

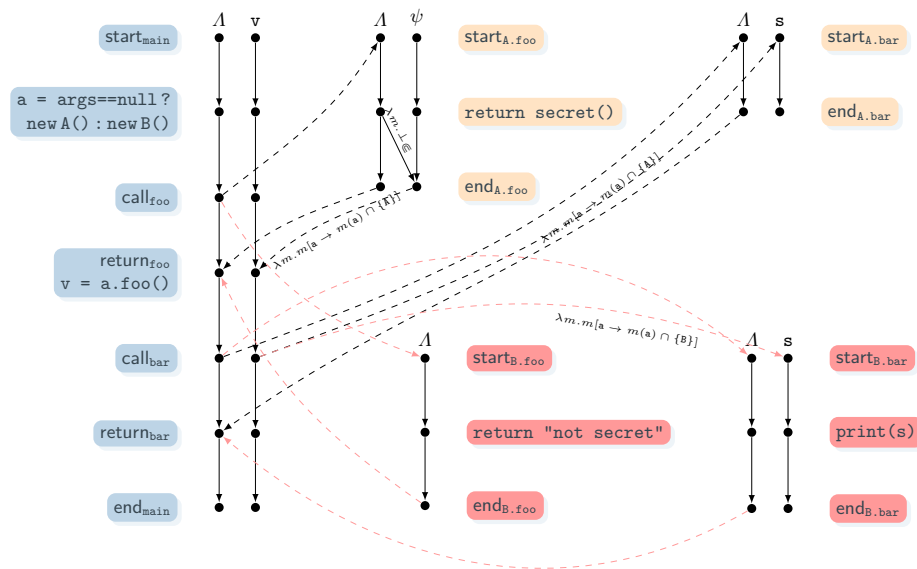
$$\tau(s_{\mathcal{F}}, f) = \{t \mid \text{lookup}(s_{\mathcal{F}}, t) = f\}. \quad (16)$$

In other words,  $\tau$  computes the set of types for which calls to methods with signatures  $s_{\mathcal{F}}$  are dispatched to  $f$ .

If there is a supergraph path from a method call with signature  $s_{\mathcal{F}}$  to the start of  $f$ , then the set  $\tau(s_{\mathcal{F}}, f)$  is always non-empty.

<sup>6</sup> See Sagiv et al. [20] for a formal definition of the representation relation for environment transformers.

I think the rest of this section will not make sense to anyone who does not already understand it. Should I try to replace it with one paragraph that tries to give an intuitive understanding of what the IDE algorithm is doing? (M)



**Fig. 3:** An example program demonstrating correlated-call edge functions on the  $\Lambda$ -node path for Listing 1. All non-labeled edges are implicitly labeled with identity functions `id`. The variable corresponding to an initial secret value is denoted as  $\psi$ .

**Definition 3.** A call site is called *monomorphic* if it can be dispatched to only one method. If a call site can be dispatched to more than one method it is called *polymorphic*.

Let  $r.c()$  be a call on a receiver  $r \in R$  with a method signature  $s_{\mathcal{F}}$  to a function  $f \in \mathcal{F}$ . If the call site is monomorphic, then  $\tau(s_{\mathcal{F}}, f)$  contains all types  $T' \subseteq T$  that are compatible with the static type of  $r$ . If the call site is polymorphic, then  $\tau(s_{\mathcal{F}}, f) \subset T'$ , since some types  $t \in T'$  cause dispatch to a method other than  $f$ .

incorporate into previous section (O)

### 3.4 Defining IFDS and IDE Problems

In Section 3, we defined what IFDS and IDE problems are, their applications, and their constraints. In this section, we describe how to create instances of IFDS and IDE problems.

**Defining an IFDS Problem** Recall that an IFDS problem instance is defined as a five-tuple

$$(G^*, D, F, M_F, \sqcap),$$

where  $G^* = (N^*, E^*)$  is the control-flow supergraph of the program,  $D$  is the set of dataflow facts,  $F \subseteq 2^D \rightarrow 2^D$  is a set of distributive dataflow functions, and the function

$$M_F : E^* \rightarrow (2^D \rightarrow 2^D)$$

maps the supergraph edges to dataflow functions, and is extended to paths by composition.

In practice, an IFDS problem can be defined by providing an exploded supergraph  $G^\# = (N^\#, E^\#)$ . Each node of  $G^\#$  is a pair  $(n, d)$ , where  $n \in N^*$  is a node in the control-flow supergraph and  $d \in (D \cup \{\mathbf{0}\})$ ,  $\mathbf{0} \notin D$ , where  $\mathbf{0}$  is an auxiliary fact that is necessary for the IFDS algorithm.

The meaning of an edge in the exploded supergraph is the following. Let  $(n_1, d_1)$  and  $(n_2, d_2)$  be two nodes in the exploded supergraph  $G^\#$ . Furthermore, assume that if fact  $d_1$  at node  $n_1$  holds, then the fact  $d_2$  at node  $n_2$  also holds. Then there is an edge  $(n_1, d_1), (n_2, d_2) \in E^\#$ .

**Defining an IDE Problem** An IDE problem instance is a four-tuple

$$(G^*, D, L, M_{\text{Env}}),$$

where  $G^*$  and  $D$  are defined in the same way as for IFDS.  $L$  is a finite-height lattice that represents the values to which dataflow facts are mapped in an IDE problem. An environment  $\text{Env}(D, L)$  maps dataflow facts to lattice elements. Finally, the map

$$M_{\text{Env}} : E^* \rightarrow (\text{Env}(D, L) \rightarrow \text{Env}(D, L))$$

is a function from the control-flow-supergraph edges to environment transformers, extended to paths by composition.

An IDE problem can be defined with a labeled exploded supergraph<sup>7</sup>, in which an edge function

$$\text{EdgeFn} : E^\# \rightarrow (L \rightarrow L) \quad (17)$$

pairs edges with *micro functions*, and is extended to paths by composition.

The set of micro functions of an IDE problem is a subset of  $L \rightarrow L$  that is closed under function meet and composition.

The meaning of an edge in the labeled exploded supergraph is the following. Let  $e = ((n_1, d_1), (n_2, d_2)) \in E^\#$  be an edge in the exploded supergraph with label  $f = \text{EdgeFn}(e)$ . Then

- if at node  $n_1$  the fact  $d_1$  was mapped to a lattice element  $l_1$  by an environment  $\text{Env}(D, L)$ , then the fact  $d_2$  at node  $n_2$  should be mapped to  $f(l_1)$ .

As shown in Sagiv et al. [20], the relationship between environment transformers and edge functions can be described with the following equations. For individual edges  $(n_1, n_2) \in E^*$ ,

$$\begin{aligned} M_{\text{Env}}((n_1, n_2))(\text{env})(d) \\ = \text{EdgeFn}((n_1, \perp), (n_2, d))(\top) \sqcap \bigcap_{d' \in D} \text{EdgeFn}((n_1, d'), (n_2, d))(\text{env}(d')), \end{aligned} \quad (18)$$

where  $\text{env}$  is an environment  $\text{Env}(D, L)$ . Informally, for a given control-flow-supergraph edge  $e$  and data-flow fact  $d$ , the  $M_{\text{Env}}$  function captures the meet of the edge function applied to all possible exploded-graph edges along  $e$ .

For paths  $p$  that start with the entry point  $\text{start}_{\text{main}}$ ,

$$M_{\text{Env}}(p)(\Omega)(d) = \bigcap_{r \in \text{RP}(p, d)} \text{EdgeFn}(r)(\top), \quad (19)$$

where  $n \in N^*$ ,  $d \in D$ ,  $p \in \text{VP}(n)$ , and  $\text{RP}$  is the set of all inter-procedurally realizable paths.

To summarize, an IDE problem can be defined by a labeled exploded supergraph

$$(G^\#, L, \text{EdgeFn}), \quad (20)$$

where each edge of the exploded supergraph corresponds to a micro function.

<sup>7</sup> The exploded supergraph in an IDE problem is defined in the same way as in an IFDS problem. The only difference is that the  $\mathbf{0}$  fact is denoted as  $\perp$  [17,20].

## 4 Correlated Calls Analysis

The correlated-calls analysis is defined as a transformation from an arbitrary IFDS problem to a corresponding IDE problem. The solution of the IDE problem is converted to a solution of the original IFDS problem. The converted IFDS result can be more precise than the original IFDS result because it avoids infeasible paths corresponding to correlated calls.

### 4.1 Transformations from IFDS to IDE

Here, you defined the output of a transformation as just  $G^\#$  and an edge function. Later, you also include the lattice. Should we include the lattice here? (O)

Let  $G^\#$  be the exploded supergraph of an IFDS problem. A *transformation*  $\mathcal{T} : (G^\#) \rightarrow (G^\#, \text{EdgeFn})$  converts it into an equivalent IDE problem. We consider two transformations:

- an equivalence transformation  $\mathcal{T}^\equiv$  (pronounced “t-equiv”) that generates IDE problems with the same precision as the original IFDS problem, and
- a correlated-call transformation  $\mathcal{T}^\subseteq$  (pronounced “t-c-c”) that generates IDE problems that exclude infeasible paths.

Both transformations keep the exploded supergraph  $G^\#$  the same, and only generate different edge functions.

**Equivalence Transformation** The lattice for the equivalence transformation  $\mathcal{T}^\equiv$  is the two-point lattice  $L^\equiv = \{\perp, \top\}$ , where  $\perp$  means “reachable”, and  $\top$  means “not reachable”. The edge functions  $\text{EdgeFn}^\equiv$  are defined as

$$\text{EdgeFn}^\equiv = \begin{cases} \lambda e. \lambda m. \perp & \text{if } d_1(e) = \Lambda \text{ and } d_2(e) \neq \Lambda; \\ \lambda e. \text{id} & \text{otherwise,} \end{cases} \quad (21)$$

where  $d_1(e)$  is the source fact of an edge  $e$  and  $d_2(e)$  is its target fact. At a “diagonal” edge from a  $\Lambda$ -fact to a non- $\Lambda$ -fact  $d$ , the micro function returns  $\perp$  to make the fact  $d$  reachable. All other micro-functions are the identity. The equivalence transformation is thus defined as:  $\mathcal{T}^\equiv((G^\#)) = (G^\#, L^\equiv, \text{EdgeFn}^\equiv)$ .

**Correlated-Calls Transformation** In the correlated-calls transformation, the lattice elements are maps from receivers to sets of types:  $L^\subseteq = \{m : R \rightarrow 2^T\}$ , where  $R$  is the set of receivers and  $T$  is the set of all types. For each receiver  $r$ , the map gives an overapproximation of the possible runtime types of  $r$ . Sets of types are ordered by the superset relation, and this is lifted to maps from receivers to sets of types, so the bottom element  $\perp_\subseteq$  maps every receiver to the set of all types, and the top element  $\top_\subseteq$  maps every receiver to the empty set of types. During an actual execution, every receiver  $r$  points to an object of some runtime type. Therefore, a data flow fact is unreachable along any feasible path if its corresponding lattice element maps any receiver to the empty set of types.

What if  $r$  points to null? (O)

A micro-function  $f \in L^\subseteq \rightarrow L^\subseteq$  defines how the map from receivers to types should be updated when an instruction is executed. The micro-function



for most kinds of instructions is the identity. On a call to and return from a specific method  $m$  called on receiver  $r$ , the micro-function restricts the receiver-to-type map to map  $r$  only to types consistent with the polymorphic dispatch to method  $m$ . Finally, when an instruction assigns an object of unknown type to a receiver  $r$ , the corresponding micro-function updates the map to map  $r$  to the set of all types. This is made precise by the following definition:

**Definition 4.** *Given a previously fixed set  $S \subseteq R$  of receivers, the micro-function  $\varepsilon_S(e)$  of a supergraph edge  $e$  is defined as:*

$$\varepsilon_S(e) = \lambda m . \begin{cases} m[r \rightarrow m(r) \cap \tau(s_{\mathcal{F}}, f)], & \text{if } e \text{ is a call-start edge, } r.c() \text{ is the call site at } n_1(e), f \text{ is the called procedure with signature } s_{\mathcal{F}}, \text{ and } r \in S; \\ m[r \rightarrow m(r) \cap \tau(s_{\mathcal{F}}, f)] & \text{if } e \text{ is an end-return edge, } r.c() \text{ is the call corresponding to the return node at } n_2(e), f \text{ is the called method with signature } s_{\mathcal{F}}, v_1, \dots, v_k \in S \text{ are the local variables in } f, \text{ and } r \in S; \\ [v_1 \rightarrow \perp_T] \dots [v_k \rightarrow \perp_T], & \\ m[r \rightarrow \perp_T], & \text{if } n_1(e) \text{ contains an assignment for } r \in S; \\ m & \text{otherwise.} \end{cases} \quad (22)$$

define  $n_1(e)$  globally somewhere early on, and consider using a more descriptive name like  $\text{src}$  instead of  $n_1$  ( $\mathbf{O}$ )

In the above definition, the purpose of the set  $S$  is to limit the set of considered receivers. We will use  $S$  in Section 4.5.

We can now define **EdgeFn**, which assigns a micro-function to each edge in the exploded supergraph. Along a  $\Lambda$ -edge, the micro function is the identity. All other functions can be described with  $\varepsilon_S$ . On a “diagonal” edge from  $\Lambda$  to a non- $\Lambda$  fact that corresponds to some data flow fact becoming reachable,  $\varepsilon_S(e)$  is applied to the initial map  $\perp_{\in}$  that conservatively allows every receiver to point to an object of any type. On all other edges,  $\varepsilon_S(e)$  is applied to the existing map before the edge. The is formalized in the following definition.

**Definition 5.** *For each edge  $e = (n_1, d_1) \rightarrow (n_2, d_2)$ ,  $\text{EdgeFn}_S^{\infty}(e)$  is defined as follows:*

$$\text{EdgeFn}_S^{\infty}(e) = \begin{cases} id & \text{if } d_1 = d_2 = \Lambda, \\ \lambda m . \varepsilon_S(e)(\perp_{\in}) & \text{if } d_1 = \Lambda \text{ and } d_2 \neq \Lambda, \\ \lambda m . \varepsilon_S(e)(m) & \text{otherwise.} \end{cases} \quad (23)$$

*Example 6.* Consider the program from Listing 1, whose exploded supergraph appeared in Figure 3.

Returning a secret value in method **A.foo** creates a “diagonal” edge from the  $\Lambda$ -fact to the secret fact  $\psi$ . The diagonal edge is labeled with the micro function  $\lambda m . \perp_{\in}$ . Thus, at the end node of the method, every receiver is mapped to the set of all types  $\perp_T$ .

On the end-return edge from **A.foo** to **main**, the set of types for the receiver **a** is restricted by the micro function  $\lambda m . m[\mathbf{a} \rightarrow m(\mathbf{a}) \cap \{\mathbf{A}\}]$  on the edge corresponding to the assignment of the return value  $\psi$  to **v**.

Similarly, on the call-start edge from `main` to `B.bar`, the possible types of the receiver `a` are further restricted by the micro-function  $\lambda m. m[\mathbf{a} \rightarrow m(\mathbf{a}) \cap \{\mathbf{B}\}]$ . on the edge that passes the argument `v` to the parameter `s`.

The composition of these micro functions results in the empty set as the possible types of the receiver `a`, indicating that this data flow path that would result in an information leak is actually infeasible.

Finally, the correlated-calls transformation is defined as  $\mathcal{T}_S^\mathbb{E}((G^\#)) = (G^\#, L_S^\mathbb{E}, \text{EdgeFn}_S^\mathbb{E})$ .

## 4.2 Converting IDE Results to IFDS Results

For each program point  $n$ , the result of an IFDS analysis gives a set of facts  $d$  that may be reached at  $n$ . The result of an IDE analysis pairs each such fact  $d$  with a lattice element  $\ell$ . Formally, for an IFDS problem  $P$ , the result  $\mathcal{R}_{\text{IFDS}}$  has type  $N \rightarrow D$ . Similarly, for an IDE problem  $Q$ , the result  $\mathcal{R}_{\text{IDE}}$  has type  $N \rightarrow D \times L$ .

Recall that in the equivalence transformation lattice  $L^\equiv$ ,  $\perp$  means reachable and  $\top$  means unreachable. Therefore, a result  $\rho$  of the equivalence IDE analysis is converted to an IFDS result as follows:  $\mathcal{U}^\equiv(\rho) = \lambda n. \{d \mid \rho(n) = (d, \top)\}$ .

In the correlated-calls transformation lattice  $L^\mathbb{E}$ , a map that maps any receiver to the empty set of possible types means that the corresponding data flow path is infeasible. Therefore, a result  $\rho$  of the correlated-calls IDE analysis is converted to an IFDS result as follows:  $\mathcal{U}^\mathbb{E}(\rho) = \lambda n. \{d \mid \rho(n) = (d, \ell), \forall r. \ell(r) \neq \top\}$ .

## 4.3 Implementation of Correlated Calls Micro-Functions

move material from following section here (O)

Conceptually, micro-functions are functions from  $L$  to  $L$ , where  $L$  is the IDE lattice, either  $L^\equiv$  or  $L^\mathbb{E}$  in our context. However, the IDE algorithm requires an efficient representation of micro-functions. The chosen representation needs to support the basic micro-functions that we presented in Section 4.1. The representation must also support function application, comparison, and be closed under function composition and meet. We now propose such a representation for the correlated-calls micro-functions.

The representation of a micro-function is a map from receivers to pairs  $\langle I, U \rangle$ , where we require  $U \subseteq I$ . The meaning of such a map  $\lambda r. \langle I_r, U_r \rangle$  is the micro-function  $\lambda m. \lambda r. (m(r) \cap I_r) \cup U_r$ . All of the micro-functions in definition 4 can be represented in this form.

Say how application, comparison, composition, and meet are implemented. Decide how much detail and how much proof to give. Also need to introduce normalization. (O)

## 4.4 Theoretical Results

In the following lemma we show that the result of an IDE problem obtained through a correlated-calls transformation is a subset of the original IFDS result.

**Lemma 1 (Precision).** *For an IFDS problem  $P$  and all  $n \in N^*$ ,*

$$\mathcal{U}^\mathbb{E}(\mathcal{R}(\mathcal{T}_R^\mathbb{E}(P)))(n) \subseteq \mathcal{R}_{\text{IFDS}}(P)(n). \quad (24)$$

Should we say  $\mathcal{R}_{\text{IDE}}$  instead of  $\mathcal{R}$  everywhere, to be consistent with  $\mathcal{R}_{\text{IFDS}}$ ? (O)

We agreed to move the proofs to the appendix. (O)

*Proof.* The transformation  $\mathcal{U}^\subseteq$  is the same as  $\mathcal{U}^\equiv$ , except that it can remove data-flow facts from the result:

$$\begin{aligned} \mathcal{U}^\subseteq(\mathcal{R}(\mathcal{T}_R^\subseteq(P)))(n) &= \{(n', D_n'^\subseteq(\mathcal{R}(\mathcal{T}_R^\subseteq(P)))) \mid n \in N^*\}(n) \\ &= D_n^\subseteq(\mathcal{R}(\mathcal{T}_R^\subseteq(P))) \\ &\subseteq \text{MVP}_F(n) \\ &= \mathcal{R}_{\text{IFDS}}(P)(n). \end{aligned} \quad \square$$

We will next show that our analysis is sound, i.e. that the result of an IDE problem obtained through a correlated-calls transformation removes only facts that occur on infeasible paths.

**Lemma 2 (Soundness).** *Let  $p = [\text{start}_{\text{main}}, \dots, n]$  be a concrete execution path, and let  $d \in D$ . If  $d \in M_F(p)(\emptyset)$ , then*

$$d \in \mathcal{U}^\subseteq(\mathcal{R}(\mathcal{T}_R^\subseteq(P)))(n). \quad (25)$$

#### 4.5 Correlated-Call Receivers

We will now show that in a correlated-calls transformation, it is enough to consider only some of the receivers of set  $R$ .

should this be a section or subsection? should it be moved somewhere else (or maybe not included at all)? (O)

**Definition 6.** *Let  $c_1$  and  $c_2$  be two call sites on a receiver  $r \in R$ . If both call sites are polymorphic, then we say that  $r$  is a correlated-call receiver.*

In other words, a correlated-call receiver is a receiver that has at least two polymorphic call invocations. We will denote the set of correlated-call receivers as  $R^\subseteq$ .

We will describe a “reduced” correlated-calls transformation in which we only consider receivers from  $R^\subseteq$  and ignore other receivers of  $R$ . We will show that IDE problems obtained through ordinary and reduced correlated-calls transformations yield the same results. In other words, we show that if a correlated calls analysis considers only correlated-call receivers, no precision is lost.

**Lemma 3.** *Let  $P$  be an IFDS problem. Then*

$$\mathcal{U}^\subseteq(\mathcal{R}(\mathcal{T}_{R^\subseteq}^\subseteq(P))) = \mathcal{U}^\subseteq(\mathcal{R}(\mathcal{T}_R^\subseteq(P))). \quad (26)$$

To summarize, Lemma 2 shows that the result  $\mathcal{R}_\subseteq$  of a correlated-calls analysis is sound since it overapproximates the data flow of all possible concrete execution paths. We have also shown in Lemma 1 that the correlated-calls analysis improves the precision of the original IFDS result  $\mathcal{R}_{\text{IFDS}}$ , because the correlated-calls result  $\mathcal{R}_\subseteq$  underapproximates an equivalence-IDE result  $\mathcal{R}_\equiv = \mathcal{R}_{\text{IFDS}}$ . Finally, we showed in Lemma 3 that a correlated-call transformation to IDE that considers only correlated-call receivers  $R^\subseteq$  achieves the same result  $\mathcal{R}_\subseteq$  that is obtained when considering all receivers  $R$ .

## 5 Correlated Calls Representations

In order to define a correlated-calls transformation, we need to represent lattice elements  $L_{R^\subseteq}^\subseteq : R^\subseteq \rightarrow 2^T$  of the target IDE problem, which are functions from receivers to sets of types, and micro functions  $L_{R^\subseteq}^\subseteq \rightarrow L_{R^\subseteq}^\subseteq$ .

As defined in Sagiv et al. [20], a representation of micro functions is efficient if the following conditions hold:

1. There is a representation for the identity and top functions.
2. The representation is closed under the meet and composition operations.
3. The micro functions form a finite-height lattice.
4. The apply, meet, composition, and equality-check operations can be computed in constant time.
5. There is a constant bound on the storage space for a micro function representation.

We will distinguish the representation of a concept from its denotation. For a concept  $c$ , we will write  $\llbracket c \rrbracket$  for its denotation and just  $c$  for its representation. For example, if we want to represent a constant function  $g$  with the constant value  $v$  that it returns, we will write for  $g$ 's representation,  $g = v$ , and for  $g$ 's denotation,  $\llbracket g \rrbracket = \lambda x. v$ .

### 5.1 Lattice Elements

Elements of the  $L_{R^\subseteq}^\subseteq$  lattice can be represented with a map from receivers to sets of types. The bottom element maps each receiver to the set  $T$  of all types:

$$\perp_\subseteq = \{(r, T) \mid r \in R^\subseteq\} \quad (27)$$

and the top element maps each receiver to the empty set:

$$\top_\subseteq = \{(r, \emptyset) \mid r \in R^\subseteq\}. \quad (28)$$

### 5.2 Micro Functions

In the context of the correlated-calls transformation to an IDE problem, a lattice element is a map from receivers to sets of types. Thus, a micro function transforms, or updates, an existing receiver-to-types map with new information about the types of a receiver.

We will represent micro functions with *update maps* which we describe in the next section.

**Update Maps** To represent micro functions that transform maps from receivers to sets of types, we use *update maps*. To define update maps, we first introduce the notions of *update functions* and *normalization*.

Let  $f$  be a micro function,  $r \in R^\subseteq$  a correlated-call receiver, and  $T$  the set of types in a program.

**Definition 7.** A non-normalized update function  $\text{update}_{f,r}^*$  is a pair of sets

$$\text{update}_{f,r}^* = \langle I_{f,r}, U_{f,r} \rangle, \quad (29)$$

where  $I_{f,r} \subseteq T$  is called the intersection set and  $U_{f,r} \subseteq T$  the union set of the update function.

**Definition 8.** Let  $\langle I, U \rangle$  be a pair of sets. The normalization function  $\mathcal{N}$  is defined as

$$\mathcal{N}(\langle I, U \rangle) = \langle I \cup U, U \rangle. \quad (30)$$

**Definition 9.** An update function  $\text{update}_{f,r}$  is a normalized pair of sets

$$\text{update}_{f,r} = \mathcal{N}(\text{update}_{f,r}^*) = \mathcal{N}(\langle I_{f,r}, U_{f,r} \rangle). \quad (31)$$

**Definition 10.** The update map of  $f$  is a map from receivers to update functions:

$$\text{update}_f = \{(r, \text{update}_{f,r}) \mid r \in R^\subseteq\}. \quad (32)$$

Thus, each micro function  $f$  is represented with an update map:

$$f = \llbracket \text{update}_f \rrbracket. \quad (33)$$

**Denotation of Update Maps** Intuitively, the meaning of a micro function is the update that it performs on a receiver-to-types map. To represent a micro function, it is enough to specify how the set of types for a given receiver has to be transformed. This is what the update map does.

For a micro function  $f$ , an update map takes a receiver and returns an update function:

$$\llbracket \text{update}_f \rrbracket = \lambda r. \llbracket \text{update}_{f,r} \rrbracket. \quad (34)$$

Given a micro function  $f : (R^\subseteq \rightarrow 2^T) \rightarrow (R^\subseteq \rightarrow 2^T)$ , an update map is defined so that

$$f(m) = \{(r, \llbracket \text{update}_f \rrbracket(r)(m(r)) \mid r \in R^\subseteq\}, \quad (35)$$

where the update map  $\llbracket \text{update}_f \rrbracket$  has type  $R^\subseteq \rightarrow (2^T \rightarrow 2^T)$ , and the update function  $\llbracket \text{update}_{f,r} \rrbracket = \llbracket \text{update}_f \rrbracket(r)$  has type  $2^T \rightarrow 2^T$ .

For any receiver-to-types map  $m$ , an update function specifies two things:

- which elements of  $m(r)$  should be preserved, and
- which new elements should be added to  $m(r)$ .

This can be achieved by maintaining the intersection  $I_{f,r}$  and union set  $U_{f,r}$ , where

$$\llbracket \text{update}_{f,r} \rrbracket(m(r)) = (m(r) \cap I_{f,r}) \cup U_{f,r}. \quad (36)$$

However, as we will see in Section 5.2, we need to be able to check update functions for equality, which is difficult to do with non-normalized update functions.

*Example 7.* For a non-empty set of types  $T$ , consider two update functions

$$u_1 = \langle T, T \rangle$$

and

$$u_2 = \langle \emptyset, T \rangle.$$

The denotations of the functions look as follows:

$$\llbracket u_1 \rrbracket = \lambda t. (t \cap T) \cup T$$

and

$$\llbracket u_2 \rrbracket = \lambda t. (t \cap \emptyset) \cup T.$$

We can see that both  $\llbracket u_1 \rrbracket$  and  $\llbracket u_2 \rrbracket$  are equal to the function  $\lambda t. T$ .

Therefore, the same function can have more than one non-normalized representation. This means that to check two functions for equality, it is not enough to compare their non-normalized representations.

This is why Definition 9 requires update functions to be normalized. Normalization makes the union set of the update function a subset of the intersection set. As we show later, normalization guarantees that each update function has a unique representation.

**Equality of Micro Functions** We will now show that micro functions can be checked for equality if their representations use normalized update functions.

First, let us show that normalization does not change the behaviour of an update function. This means that the normalized and non-normalized versions of an update function always denote the same function.

**Lemma 4.**  $\llbracket \mathcal{N}(\text{update}_{f,r}^*) \rrbracket = \llbracket \text{update}_{f,r}^* \rrbracket.$

Let us show that two update functions are equal if and only if their normalized representations are equal.

It is obvious that two functions represented with the same pairs of sets are equal. The following lemma states that two different pairs of sets represent different update functions.

**Lemma 5.** *Let  $\langle I, U \rangle$  and  $\langle I', U' \rangle$  be two normalized update functions such that  $\langle I, U \rangle \neq \langle I', U' \rangle$ . Then  $\llbracket \langle I, U \rangle \rrbracket \neq \llbracket \langle I', U' \rangle \rrbracket$ .*

We have shown that transfer functions can be represented using update maps.

**Operations on Update Maps** Let us now define the apply, compose, meet, and equals functions on micro function representations.

For an update map  $f$  and a receiver  $r \in R^\subseteq$ , let the update function  $f(r) = \langle I_{f,r}, U_{f,r} \rangle$ .

Then the operations on the update maps  $f_1$  and  $f_2$  are defined as follows:

$$\text{apply}_{f_1} = \lambda m. \{ (r, (m(r) \cap I_{f,r}) \cup U_{f,r}) \mid r \in R^\infty \}, \quad (37)$$

$$f_1 \circ f_2 = \{ (r, \mathcal{N}(\langle I_{f_1,r} \cap I_{f_2,r}, (I_{f_1,r} \cap U_{f_2,r}) \cup U_{f_1,r} \rangle)) \mid r \in R^\infty \}, \quad (38)$$

$$f_1 \sqcap f_2 = \{ (r, \langle I_{f_1,r} \cup I_{f_2,r}, U_{f_1,r} \cup U_{f_2,r} \rangle) \mid r \in R^\infty \}, \quad (39)$$

$$\text{equals}(f_1, f_2) = \begin{cases} \text{true} & \text{if } f_1 \text{ and } f_2 \text{ are structurally equal,} \\ \text{false} & \text{otherwise.} \end{cases} \quad (40)$$

The denotation of the operations on update maps can be explained in the following way.

The apply function of a micro function  $f$  maps over all receivers. For each receiver  $r \in R^\infty$ ,  $\text{update}_f(r)$  transforms the argument receiver-to-types map  $m$ . It returns a new map in which  $r$  is mapped to a new set of types,  $\text{update}_{f,r}(m(r))$ :

$$\begin{aligned} \llbracket \text{apply}_f \rrbracket &= \llbracket \lambda m. \{ (r, (m(r) \cap I_{f,r}) \cup U_{f,r}) \mid r \in R^\infty \} \rrbracket \\ &= \lambda m. \{ (r, \llbracket \text{update}_{f,r} \rrbracket (m(r))) \mid r \in R^\infty \}. \end{aligned}$$

Note that in the beginning of the algorithm,  $m$  maps each receiver to  $\perp_T$  (all types).

Composing two micro functions means to compose their update-map denotations:

$$\begin{aligned} \llbracket f_1 \circ f_2 \rrbracket &= \llbracket \{ (r, \mathcal{N}(\langle I_{f_1,r} \cap I_{f_2,r}, (I_{f_1,r} \cap U_{f_2,r}) \cup U_{f_1,r} \rangle)) \mid r \in R^\infty \} \rrbracket \\ &= {}^8 \lambda m. \{ (r, \llbracket \langle I_{f_1,r} \cap I_{f_2,r}, (I_{f_1,r} \cap U_{f_2,r}) \cup U_{f_1,r} \rangle \rrbracket (m(r))) \mid r \in R^\infty \} \\ &= \lambda m. \{ (r, (I_{f_1,r} \cap m(r) \cap I_{f_2,r}) \cup (I_{f_1,r} \cap U_{f_2,r}) \cup U_{f_1,r}) \mid r \in R^\infty \} \\ &= \lambda m. \{ (r, ((m(r) \cap I_{f_2,r}) \cup U_{f_2,r}) \cap I_{f_1,r} \cup U_{f_1,r}) \mid r \in R^\infty \} \\ &= (\lambda m. \{ (r, (m(r) \cap I_{f_1,r}) \cup U_{f_1,r}) \mid r \in R^\infty \}) \\ &\quad \circ (\lambda m. \{ (r, (m(r) \cap I_{f_2,r}) \cup U_{f_2,r}) \mid r \in R^\infty \}) \\ &= \llbracket \{ (r, \langle I_{f_1,r}, I_{f_2,r} \rangle) \} \rrbracket \circ \llbracket \{ (r, \langle I_{f_2,r}, U_{f_2,r} \rangle) \} \rrbracket \\ &= \llbracket f_1 \rrbracket \circ \llbracket f_2 \rrbracket. \end{aligned} \quad (41)$$

Should this be somehow moved to the appendix? Or should we just say "It can be shown that ..." and remove the whole computations at all? (M)

Finally, the meet operation on two micro functions is the union of their update maps:

$$\begin{aligned}
\llbracket f_1 \sqcap f_2 \rrbracket &= \llbracket \{ (r, \langle I_{f_1, r} \cup I_{f_2, r}, U_{f_1, r} \cup U_{f_2, r} \rangle) \mid r \in R^\infty \} \rrbracket \\
&= \lambda m. \{ (r, m(r) \cap (I_{f_1, r} \cup I_{f_2, r}) \cup U_{f_1, r} \cup U_{f_2, r}) \mid r \in R^\infty \} \\
&= \lambda m. \{ (r, ((m(r) \cap I_{f_1, r}) \cup U_{f_1, r}) \cup ((m(r) \cap I_{f_2, r}) \cup U_{f_2, r})) \mid r \in R^\infty \} \\
&= \lambda m. \{ (r, (m(r) \cap I_{f_1, r}) \cup U_{f_1, r}) \mid r \in R^\infty \} \\
&\quad \sqcap \lambda m. \{ (r, (m(r) \cap I_{f_2, r}) \cup U_{f_2, r}) \mid r \in R^\infty \} \\
&= \llbracket f_1 \rrbracket \sqcap \llbracket f_2 \rrbracket.
\end{aligned} \tag{42}$$

We can now show how the correlated-calls definitions of the meet and composition operations on micro functions allow us to detect infeasible paths in a program.

*Example 8.* The edges of the exploded supergraph in Figure 3 correspond to the edges of an IFDS taint analysis. We can see that there is a path from the node (`startmain`,  $\Delta$ ) to (`print(s)`,  $s$ ). This means that the IFDS taint analysis considers  $s$  to be a secret value that is leaked at the print statement.

The correlated-calls analysis, on the other hand, detects that the path to (`print(s)`,  $s$ ) is infeasible: at the print node, the lattice element corresponding to the fact  $s$  contains a mapping  $a \rightarrow \top_T$ .

The lattice element for the print statement is evaluated as follows:

$$\begin{aligned}
&((\lambda m. m[a \rightarrow m(a) \cap \{B\}]) \circ \text{id} \circ (\lambda m. m[a \rightarrow m(a) \cap \{A\}])) \\
&\quad \circ (\lambda m. \perp_\infty) \circ \text{id} \circ \dots \circ \text{id})(\top_\infty) \\
&= (\perp_\infty[a \rightarrow m(a) \cap (\{A\} \cap \{B\})]) \\
&= (\perp_\infty[a \rightarrow \top_T]).
\end{aligned}$$

Therefore, the path to the print statement will be considered infeasible, and the analysis does not claim that the program leaks a secret value.

The following Lemma states that the representation of micro functions is efficient according to the definition of efficiency discussed in Section 5.

**Lemma 6.** *The correlated-call representation of a micro function is efficient.*

*Final Remarks* A straightforward solution to representing micro functions would be to use the function constructs that are provided by many programming languages. The efficiency requirement prohibits us from doing so. For most programming languages, equality for functions is either defined as reference equality (as in Scala), or is not defined at all (as in Haskell). Even if we were to define our own definition of equality for functions, we would have to iterate over the whole domain of the functions and compare the results of the function applications, which

<sup>8</sup> See Lemma 4.



would be inefficient. Additionally, the equality check would be non-terminating if the domain of the functions were infinite, and undecidable if the language for defining the functions were Turing-complete.

Second, a composition  $f$  of two functions  $f_1$  and  $f_2$  would have to store both  $f_1$  and  $f_2$ . For instance, if  $f_1 = \lambda x . x + 1$  and  $f_2 = \lambda x . x + 2$ , then  $f = f_2 \circ f_1$  would be represented as

$$f = \lambda x . (\lambda y . y + 2)((\lambda z . z + 1) x)$$

instead of

$$f = \lambda x . x + 3.$$

Having a compact representation for function composition is especially important for the first phase of the IDE algorithm, in the computation of jump functions [20]. The same argument applies to computing function meets.

**Edge Function Representation** We will now show the representations for the correlated-call micro functions  $\text{EdgeFn}_{R^\infty}^\infty(e)$ , described in Definition 5. Let  $\text{identity} = \langle \perp_T, \top_T \rangle$  represent the identity function  $\text{id}$  and  $\text{bottom} = \langle \perp_T, \perp_T \rangle$  represent the function  $\lambda t . \perp_T$ .

On the call-start edge,

$$\begin{aligned} & m[r \rightarrow (m(r) \cap \tau(s_{\mathcal{F}}, f))] \\ &= \llbracket \{(r, \langle \tau(s_{\mathcal{F}}, f), \top_T \rangle)\} \cup \{(r', \text{identity}) \mid r' \in R^\infty, r' \neq r\} \rrbracket. \end{aligned} \quad (43)$$

On the end-return edge,

$$\begin{aligned} & \lambda m . m[v_1 \rightarrow \perp_T] \dots [v_k \rightarrow \perp_T][r \rightarrow (m(r) \cap \tau(s_{\mathcal{F}}, f))] \\ &= \llbracket \{(r, \langle \tau(s_{\mathcal{F}}, f), \top_T \rangle)\} \cup \{(r', w(r')) \mid r \in R^\infty, r' \neq r\} \rrbracket, \end{aligned} \quad (44)$$

where

$$w(r) = \begin{cases} \text{bottom} & \text{if } r \text{ is a local variable in the exiting method,} \\ \text{identity} & \text{otherwise.} \end{cases}$$

For assignments in the source node of  $e$ ,

$$\lambda m . m[r \rightarrow \perp_T] = \llbracket \{(r, \text{bottom})\} \cup \{(r', \text{identity}) \mid r' \in R^\infty, r' \neq r\} \rrbracket. \quad (45)$$

In the default case,

$$\text{id} = \llbracket \{(r, \text{identity}), r \in R^\infty\} \rrbracket.$$

We have shown how IDE problems that account for correlated calls can be represented in an efficient way. In the next section, we address the implementation and present an evaluation of the correlated-calls analysis.

## 6 Evaluation

This section discusses implementation aspects of the correlated-calls analysis and presents experimental results.

### 6.1 Implementation of the Analysis

The correlated-calls analysis was implemented in the Scala programming language [16]. We chose Java as the target language for client programs of the analysis. To retrieve information about an input program, such as its control-flow supergraph or the set of receivers and their types, we used the WALA framework for static analysis on Java bytecode [6].

Since WALA currently only contains an implementation of IFDS, we implemented IDE from scratch. Instead of using WALA’s IFDS implementation, to run an IFDS problem, we converted it to an IDE problem and used our own IDE solver.

**IFDS** As described in Section 3.4, an IFDS problem is defined in terms of an exploded supergraph. The control-flow supergraph of an input program can be retrieved using WALA. Hence, our implementation of an IFDS problem should be able to convert a control-flow supergraph into an exploded supergraph.

We represent an IFDS problem with a trait, or protocol, that contains declarations of four *flow functions*. Each function has type

$$F : (N \times D \times N) \rightarrow 2^D$$

and defines a set of edges on the exploded graph. Given an edge  $(n_1, n_2)$  of the control-flow supergraph and the fact  $d_1$  that corresponds to the source node  $n_1$ ,  $F(n_1, d_1, n_2)$  returns the set of all facts  $d_2 \in D_2$  such that  $((n_1, d_1), (n_2, d_2)) \in E^\#$ <sup>9</sup>. The four functions are:

- **call-start**, for inter-procedural edges from a call node to the start node of the target method;
- **call-return**, for intra-procedural edges from a call node to its return node;
- **end-return**, for inter-procedural edges from the end node of a method to the return node of the callee;
- **default**, for all other intra-procedural edges.

*Taint Analysis* Using this representation of an IFDS problem, we implemented an IFDS problem instance for taint analysis. We used it as a sample IFDS problem on which to evaluate the correlated-calls-IDE construction.

<sup>9</sup> In each invocation of a flow function, the fact  $d_1$  is provided by the IDE algorithm.

We should incorporate the four functions into the explanation of IFDS in the background. Then we don't need them here. (O)

Let  $N^*$  be the control-flow supergraph of a program and  $D$  the set of the program variables. Let  $\text{encl}(n)$  be a function that returns the enclosing method of a node  $n \in N^*$ . Finally, let the function  $r_m : D \rightarrow 2^D$  be defined as follows:

$$r_m(d) = \begin{cases} \emptyset & \text{if } d \text{ is a local variable in method } m, \\ \{d\} & \text{otherwise.} \end{cases} \quad (46)$$

When defining the flow functions for a taint analysis, we will use  $r_m$  to avoid the propagation of local variables, as shown below.

For a fact  $d_1 \in D \cup \{\mathbf{0}\}$  and two nodes  $n_1, n_2 \in N^*$ , the simplified<sup>10</sup> version of flow functions for a taint-analysis looks as follows.

If  $n_1$  is a call node that calls method  $m$ , and  $n_2$  is  $m$ 's start node,

$$\text{call-start}(n_1, d_1, n_2) = \begin{cases} r_{\text{encl}(n_1)}(d_1) \cup \{v\} & \text{if } a \text{ is the } i\text{th argument of the call,} \\ & d_1 = a, \text{ and } v \text{ is the } i\text{th parameter} \\ & \text{of } m; \\ r_{\text{encl}(n_1)}(d_1) & \text{otherwise.} \end{cases}$$

If  $n_1$  is a call node with corresponding return node  $n_2$ ,

$$\text{call-return}(n_1, d_1, n_2) = \begin{cases} \{d_1\} & \text{if } d_1 \text{ is a local variable in } \text{encl}(n_1), \\ \emptyset & \text{otherwise.} \end{cases}$$

If  $c$  is a call node calling method  $m$ ,  $n_1$  is  $m$ 's end node, and  $n_2$  is  $c$ 's return node,

$$\text{end-return}(n_1, d_1, n_2) = \begin{cases} r_{\text{encl}(n_1)}(d_1) \cup \{x\} & \text{if } n_1 \text{ is a return statement} \\ & \text{returning } v, n_2 \text{ is an assignment} \\ & \text{with left-hand side } x, \text{ and } d_1 = v; \\ r_{\text{encl}(n_1)}(d_1) & \text{otherwise.} \end{cases}$$

Otherwise,

$$\text{default}(n_1, d_1, n_2) = \{d_1\}.$$

*Example 9.* Consider the supergraph in Figure ?? . The call-to-start flow function from method `main` to `f` looks as follows:

$$\begin{aligned} \text{call-start}(\text{call}_{\mathbf{a}, \mathbf{f}}, \mathbf{a}, \text{start}_{\mathbf{f}}) &= r_{\text{main}}(\mathbf{a}) \cup \{\mathbf{s}\} \\ &= \{\mathbf{s}\}. \end{aligned}$$

We can see that correspondingly, the exploded supergraph contains an edge from  $(\text{call}_{\mathbf{a}, \mathbf{f}}, \mathbf{a})$  to  $(\text{start}_{\mathbf{f}}, \mathbf{s})$ .

<sup>10</sup> For simplicity, the shown flow functions do not account for different Java-specific features such as arrays, fields, operations on strings, etc.

**IDE** The correlated-calls analysis was implemented as an IDE problem instance.

We defined an IDE problem in the same way as an IFDS problem, except that the IDE flow functions are of type

$$(N \times D \times N) \rightarrow 2^{D \times (L \rightarrow L)}.$$

With the new flow functions, we can implement a labeled exploded supergraph, since the new flow functions return a set of facts that are paired with micro functions.

For example, if  $Q$  is an IDE problem, then the call-to-start flow function for  $Q$  is defined as follows:

$$\begin{aligned} & \text{call-start}^Q(n_1, d_1, n_2) \\ &= \left\{ (d_2, f) \mid d_2 \in D, f \in L^Q \rightarrow L^Q : \text{EdgeFn}^Q((n_1, d_1), (n_2, d_2)) = f \right\}. \end{aligned}$$

The other flow functions are defined analogously.

## 6.2 Testing

In this section we assess the correctness and effectiveness of the correlated-calls analysis.

**Conversion from IFDS to IDE** We implemented the equivalence transformation  $\mathcal{T}^\equiv$  and the correlated-calls transformation  $\mathcal{T}_{R^\epsilon}^\equiv$  from IFDS to IDE described in Section 4.1. To run an IFDS problem, we converted it to an IDE problem using  $\mathcal{T}^\equiv$  and  $\mathcal{T}_{R^\epsilon}^\equiv$  and used our IDE analysis algorithm to run the latter.

Given an IFDS problem described with IFDS flow functions, an equivalence transformation creates an IDE problem described with the following IDE flow functions:

$$\begin{aligned} \text{call-start}^\equiv(n_1, d_1, n_2) &= \{(d_2, \epsilon(d_1, d_2)) \mid d_2 \in \text{call-start}(n_1, d_1, n_2)\} \\ \text{call-return}^\equiv(n_1, d_1, n_2) &= \{(d_2, \epsilon(d_1, d_2)) \mid d_2 \in \text{call-return}(n_1, d_1, n_2)\} \\ \text{end-return}^\equiv(n_1, d_1, n_2) &= \{(d_2, \epsilon(d_1, d_2)) \mid d_2 \in \text{end-return}(n_1, d_1, n_2)\} \\ \text{default}^\equiv(n_1, d_1, n_2) &= \{(d_2, \epsilon(d_1, d_2)) \mid d_2 \in \text{default}(n_1, d_1, n_2)\}, \end{aligned}$$

where  $\epsilon$  is the bottom function on an edge from a  $\Lambda$ -fact to a non- $\Lambda$ -fact, and the identity function otherwise:

$$\epsilon(d_1, d_2) = \begin{cases} \lambda l. \perp & \text{if } d_1 = \Lambda \text{ and } d_2 \neq \Lambda; \\ \text{id} & \text{otherwise.} \end{cases}$$

We also implemented a correlated-call transformation from IFDS into IDE problems that consider correlated calls. This transformation is described in Section 4.1.

**Regression Testing** We used regression tests to assess the correctness of the implemented analyses. Each test involves running a certain analysis on one input Java program.

*IDE-Implementation Correctness* To test the correctness of the IDE algorithm implementation, we implemented a copy-constant-propagation IDE problem [20]. In a copy-constant propagation analysis, a variable is considered constant if it is assigned a constant literal or another variable that is also a constant. For example, in a program

```
int a = 1;
int b = a;
int c = a + b;
int d = a + 2;
```

`a` and `b` are considered constant, but `c` and `d` are not (although `d` would be considered constant in linear-constant propagation).

We tested the propagation of constants on different intra- and inter-procedural data-flow paths, in parameter passing, and in conditional branches. Each regression test contained assertions of the form “at the end of method *m*, variable with name *x* should be (not) constant”.

We also tested the implementation of the IDE algorithm on an IDE problem generated by conversion from an IFDS problem.

To do that, we implemented an IFDS instance for taint analysis.

Recall from Section 3.2 that taint analysis aims to discover variables that are secret at a given program point called a sink.

We used assertions of the form “at program statement *n*, variable *x* should be (not) secret” by defining the sink of a secret value through special `isSecret` and `notSecret` methods. Those methods asserted that the parameter passed to them is secret and not secret, respectively. To define a source secret value we created a static `secret()` method that returned a string.

*Example 10.* Listing 4 illustrates the use of the `isSecret` and `notSecret` assertions.

We tested data flow through

- method calls and returns;
- conditional branches and loops, including nested constructions, the ternary operator, and `switch` statements;
- arrays and fields<sup>11</sup>;
- static and instance class members;

<sup>11</sup> In Java, arrays are allocated on the heap, and array elements can be aliases of each other. Hence, if any array element gets assigned a secret value, we considered all elements of any `String` or `Object` array in the program secret. For the same reason, if a field `f` of an object of class `A` is assigned a secret value, then we considered the field `f` of any object of class `A` secret.

```

public static void main(String[] args) {
    String n = "not_secret";
    notSecret(n); // assert that n is not secret
    String s1 = f(n);
    isSecret(s1);
    String s2 = f(secret());
}

static String f(String str) {
    isSecret(str);
    return str;
}

public static String secret() { // the secret source
    return "secret";
}

```

**Fig. 4:** Example usage of `isSecret` and `notSecret` assertions in regression tests

- classes and interfaces that involve inheritance, overriding, and overloading;
- recursion;
- library calls<sup>12</sup>;
- string concatenation and usage of the `StringBuffer` and `StringBuilder` classes<sup>13</sup>;
- generics, type conversions through castings, and exception handling.

Our taint analysis implementation becomes unsound in the presence of static initializers. If a static field is initialized to a secret value, our analysis will not detect it as such.

A static initializer is invoked only once, before the instance creation of a class or the access of a static member of that class. Static initializers are invoked lazily by the Java Virtual Machine [11]. This makes finding out at which program point a static initializer is invoked undecidable [7]. To account for static initializers in the analysis would require modifying WALA’s control-flow supergraph (which does not have edges to static initializers) or using a data-flow analysis for static initialization. Since the primary purpose of the taint-analysis implementation was to test the correlated-call analysis, we did not include a static-initializer analysis in this work.

<sup>12</sup> We created a specification for library functions that allowed us to indicate under which conditions a library function returned a secret value. This let us avoid the expensive analysis of library functions.

<sup>13</sup> Using mutation, objects of these classes can be converted into wrappers around secret strings. This is why we added a special handling for `StringBuffer` and `StringBuilder` objects. For instance, if a field had the `StringBuilder` type, it was considered secret.

*Correlated-Calls-Analysis Correctness* We tested the implementation of the correlated-calls analysis by converting the taint analysis into an IDE problem with an implementation of  $\mathcal{T}_{R^\infty}$ .

Since none of the test cases in the previous section contained correlated calls, we used the same tests with the same assertions to ensure that the correlated-calls analysis produces the same results as an IFDS-equivalent analysis in the absence of correlated calls.

We then added test cases that contained correlated calls. We added a new assertion method, `notSecretCC`. For the IFDS-equivalent analysis, the method asserted that the argument passed to it was secret, and for the correlated-calls analysis, it asserted that the argument was not secret.

Separately, we used unit tests to check the implementation correctness of micro functions. We wrote assertions for the results of the equality, meet, and composition operations on all possible combinations of the identity, top, bottom, and constant functions.

**Benchmark Testing** To assess the benefit of the correlated-calls analysis, we counted the frequencies of correlated-call occurrences in the Dacapo benchmarks [2]. We then ran the normal- and correlated-call-taint analysis on the Dacapo benchmarks to see what improvement we would get from the correlated-calls analysis.

*Occurrences of Correlated Calls* Our goal was to obtain an upper bound on the number of redundant IFDS-result nodes that could be potentially removed by our analysis. We counted the number of correlated calls that occurred in programs of the Dacapo benchmarks, as shown in Table 1.

In the table, the number of all call sites in a program is denoted as  $C$ . Polymorphic call sites are denoted as  $C_P$ , and correlated call sites as  $C^\infty$ . The first four columns indicate the overall number of various call sites and correlated-call receivers in a program. The last three columns indicate the ratio of polymorphic to all call sites, the ratio of correlated to polymorphic call sites, and the ratio of correlated call sites to correlated-call receivers.

We can see that on average, 3% of all call sites  $C$  are polymorphic call sites  $C_P$ . Out of those call sites, 38% are correlated call sites  $C^\infty$ . We also see that for one correlated-call receiver, there are on average three correlated calls.

*Experiments* We ran the analysis on the Dacapo benchmarks to test if the taint analysis would benefit from the improved, correlated-calls based, analysis. We defined any user input string to be considered a secret source and compared the overall number of results in the original and correlated-call taint analyses. If the number of secret values in the original result were larger than in the correlated-call result, we would see a practical benefit from our analysis.

However, even when we considered each program point as a sink, the “improved” analysis revealed the same number of secret values as the original taint analysis.

We need to decide what exactly we want to include in the evaluation. The numbers of correlated calls are interesting. But I think we want to leave out all mention of the taint analysis. Perhaps we should have no empirical results at all (and only the theory/proofs)? (O)

**Table 1:** Frequencies of correlated-call occurrences in the Dacapo benchmarks

| Benchmark         | $ C $        | $ C_P $    | $ C^\infty $ | $ R^\infty $ | $\frac{ C_P }{ C }$ | $\frac{ C^\infty }{ C_P }$ | $\frac{ C^\infty }{ R^\infty }$ |
|-------------------|--------------|------------|--------------|--------------|---------------------|----------------------------|---------------------------------|
|                   |              |            |              |              |                     |                            |                                 |
| <b>antlr</b>      | 7,610        | 428        | 299          | 70           | <b>6%</b>           | <b>70%</b>                 | 4                               |
| <b>bloat</b>      | 18,157       | 933        | 429          | 119          | <b>5%</b>           | <b>46%</b>                 | 4                               |
| <b>chart</b>      | 18,101       | 466        | 195          | 61           | <b>3%</b>           | <b>42%</b>                 | 3                               |
| <b>eclipse</b>    | 3,222        | 100        | 35           | 10           | <b>3%</b>           | <b>35%</b>                 | 4                               |
| <b>fop</b>        | 4,831        | 129        | 40           | 12           | <b>3%</b>           | <b>31%</b>                 | 3                               |
| <b>hsqldb</b>     | 3,573        | 81         | 35           | 10           | <b>2%</b>           | <b>43%</b>                 | 4                               |
| <b>kython</b>     | 12,149       | 487        | 129          | 54           | <b>4%</b>           | <b>26%</b>                 | 2                               |
| <b>luindex</b>    | 7,190        | 188        | 79           | 29           | <b>3%</b>           | <b>42%</b>                 | 3                               |
| <b>lusearch</b>   | 9,043        | 350        | 126          | 47           | <b>4%</b>           | <b>36%</b>                 | 3                               |
| <b>pmd</b>        | 10,972       | 219        | 68           | 23           | <b>2%</b>           | <b>31%</b>                 | 3                               |
| <b>xalan</b>      | 3,889        | 110        | 35           | 10           | <b>3%</b>           | <b>32%</b>                 | 4                               |
| <b>Geom. mean</b> | <b>7,572</b> | <b>240</b> | <b>91</b>    | <b>29</b>    | <b>3%</b>           | <b>38%</b>                 | <b>3</b>                        |

A correlated call that could affect a taint-analysis result could most likely occur in the following scenario:

- there is a receiver with at least two polymorphic calls;
- at least one of the calls  $c_1$  returns a string — this would mean that the method potentially returns a secret value;
- at least one of the calls  $c_2$  takes a string parameter — this would mean that a secret value could potentially be propagated to the method as an argument.

Then, if the correlated call occurred on an invocation  $c_2(c_1())$ , there might be a possibility of benefiting from the correlated-calls analysis. Given the relatively rare occurrence of correlated calls, this situation is not likely to appear often. This is illustrated in Table 2 which shows how often correlated calls would invoke methods that either take a string as a parameter *or* return a string. The set of receivers on which there are invocations of such methods is denoted as  $R^\infty_S$ . A situation where one correlated call returned a string, *and* another correlated call on the same receiver took a string parameter, appeared in only one case in the **antlr** benchmark. However, the strings invoked were not designated as secret.

This explains why, specifically for a taint analysis as the client analysis, and specifically for the Dacapo benchmarks, the correlated call analysis did not make a difference.

### 6.3 Future Work

In this section we point out the limitations of the correlated-calls analysis and suggest improvements to the analysis for future work.

One limitation of the analysis is that it only works for IFDS problems like taint analysis, reachable definitions, or available expressions. The correlated-call

We should convert this section into two or three sentences to be added to the conclusion. They should focus especially on the interprocedurally-correlated receivers. (O)



**Table 2:** Frequency of correlated-call receivers for which at least one of the correlated calls takes a string as a parameter or returns a string

| Benchmark         | $ R^{\subseteq}_S $ | $ R^{\subseteq} $ | $\frac{ R^{\subseteq}_S }{ R^{\subseteq} }$ |
|-------------------|---------------------|-------------------|---|
| <b>antlr</b>      | 43                  | 70                | 62%   |
| <b>bloat</b>      | 0                   | 119               | 0%  |
| <b>chart</b>      | 1                   | 61                | 2%  |
| <b>eclipse</b>    | 0                   | 10                | 0%  |
| <b>fop</b>        | 0                   | 12                | 0%  |
| <b>hsqldb</b>     | 0                   | 10                | 0%  |
| <b>jython</b>     | 6                   | 54                | 23%   |
| <b>luindex</b>    | 0                   | 29                | 0%  |
| <b>lusearch</b>   | 2                   | 47                | 6%  |
| <b>pmd</b>        | 1                   | 23                | 3%  |
| <b>xalan</b>      | 0                   | 10                | 0%  |
| <b>Geom. mean</b> | <b>3</b>            | <b>29</b>         | <b>9</b>                                    |

analysis is not applicable to IDE problems like copy- or linear-constant propagation. Therefore, a possible direction for future work is to create a correlated-calls analysis that transforms an original IDE problem into one that considers correlated calls (with a modified lattice and edge function definition), and then transforms the correlated-calls result into a more precise result of the original IDE problem.

Another constraint of the algorithm is that it only accounts for intra-procedurally-correlated receivers, or receivers on which correlated calls occur within one method. For example, in Listing 5, **a** is a correlated-call receiver, since there are two polymorphic method invocations on **a**. However, the first one, **a.setString()**, is inside method **main**, and the second one, **a.printString()**, is inside method **propagate**. Therefore, we do not treat **a** as a correlated-call receiver, and the analysis would not improve the original IFDS result for this program.

Finally, correlated calls can occur on multiple receivers and other scenarios discussed in [?] that are not handled in this work.

## 7 Related Work

IFDS is a version of the functional approach to data-flow analysis developed by M. Sharir and A. Pnueli [21]. Their algorithm is based on computing *summary functions* that return the data-flow value at the end of a procedure, given the data-flow value at the start of the procedure. IFDS problems form a more restricted set of data-flow problems: unlike in the functional approach, IFDS flow functions have to be distributive, and the set of data-flow facts  $D$  has to be finite. However, the IFDS algorithm is more general than Sharir's and Pnueli's algo-

should we move Related Work to near the end of the paper? (O)

```
class A {  
    String string;  
  
    public static void main(String[] args) {  
        A a = args == null ? new A() : new B();  
        a.setString();  
        propagate(a);  
    }  
  
    static void propagate(A a) {  
        a.printString();  
    }  
  
    void setString() {  
        string = secret();  
    }  
  
    void printString() {  
        System.out.println("not_secret");  
    }  
}  
  
class B extends A {  
    void setString() {  
        string = "not_secret";  
    }  
  
    void printString() {  
        System.out.println(a);  
    }  
}
```

**Fig. 5:** Inter-procedurally-correlated calls

rithm in that it can handle programs containing local variables and parameters in recursive methods.

IFDS has been used to encode a variety of data-flow problems. More complex examples of applications include typestate analysis (determining which operations can be performed on an object at a given program point) [13] or shape analysis (detecting errors and validating properties of programs at compile time) [10].

IFDS is implemented for two popular static-analysis frameworks, the T.J. Watson Libraries for Analysis (WALA) [6] and Soot [25].

WALA is a framework for static analysis on Java bytecode developed by the IBM T.J. Watson Research Center. In the implementation of our work, we use WALA to build and traverse the supergraph (a special kind of control-flow graph) of a Java program<sup>14</sup>.

Soot is a framework for program analysis and optimization on Java bytecode, developed by the Sable Research Group at McGill University. Unlike WALA, Soot also has an implementation of the IDE algorithm. The IFDS and IDE implementations for Soot are part of the Heros project [3].

Whereas one advantage of Soot’s IFDS implementation (and other static analysis tools) is ease of use and extensibility, WALA’s primary focus is efficiency. For example, WALA uses bit-vectors to represent some of the analysis data types, like local variables and parameters. Another difference is that WALA’s intermediate representation of a program uses static single assignment (SSA) form [5]. SSA form is a representation of the program in which each variable has only one definition (assignment). SSA can make dataflow analysis simpler and more efficient [1].

Work on improving the IFDS algorithm includes Practical Extensions by N. Naeem and O. Lhoták [14]. Their paper presents four extensions to the IFDS algorithm. Two of the extensions improve the efficiency of the IFDS analysis for certain classes of IFDS problems. Another extension widens the class of problems applicable for the IFDS analysis. However, those extensions do not affect the precision of IFDS problems. Our analysis, in contrast, does not improve the efficiency or generality of IFDS, but it allows us to solve IFDS problems more precisely.

The fourth extension is targeted towards programs that are represented in SSA form. Executing the IFDS analysis on such programs results in loss of precision in the presence of control-flow constructs (e.g. conditionals and loops), compared to programs in non-SSA form. The extension makes the IFDS analysis on programs in SSA form as precise as on programs that are not represented in SSA form. In contrast, the correlated-calls analysis is applicable to programs in both SSA and non-SSA forms. Even if applied to a program in SSA form, our analysis and the extension improve the precision of IFDS in unrelated situations: the first analysis handles correlated calls, and the latter handles control-flow constructs. Thus, an IFDS analysis could benefit from both precision improvements independently.

Another work on improving the efficiency of the IFDS algorithm is E. Bodden et al.’s framework for the analysis of software products lines [4]. Their paper

<sup>14</sup> However, we do not use WALA’s IFDS implementation, as explained in Section 6.

uses transformations from IFDS to IDE problems, a technique we also employ. Finally, J. Rodriguez and O. Lhoták implemented a concurrent version of the IFDS algorithm using actors [19]. However, neither of those works is concerned with improving the precision of IFDS results.

The correlated-calls analysis improves the precision of a data-flow analysis by eliminating a special type of infeasible paths. This is similar to the idea of context-sensitive analysis: just as a context-sensitive analysis eliminates infeasible paths from the end of a procedure to the call sites that do not match the given procedure call, the correlated-calls analysis eliminates infeasible paths caused by correlated method calls.

The idea of using correlated calls to remove infeasible paths in data-flow analyses of object-oriented programs was introduced by F. Tip [?]. The possibility of using IDE to achieve this is mentioned, but not elaborated upon. Our work presents a concrete solution to the problem and an implementation of that solution.

The idea of eliminating infeasible paths caused by correlated calls is similar to M. Sridharan et al.’s work on improving the precision of pointer analysis for JavaScript programs [22]. For each pointer, a pointer analysis determines the possible set of objects (the *points-to* set) that the pointer can reference at a given program point. In JavaScript, it is challenging to compute the points-to set of fields because in general, field names can be derived from arbitrary expressions and bound at runtime. As a result, an imprecise data-flow analysis will include infeasible paths between values of the form  $o[p]$  (access of a property  $p$  of object  $o$ ), where at compile time,  $p$  can be bound to different values. The idea of the paper is to track all dynamic property accesses (reads and writes) on an object  $o$  with property name  $p$ . The code snippets containing the references  $o[p]$  are then extracted into a separate function  $f$ . The analysis is then run so that for each possible value of  $p$ ,  $f$  is analyzed separately; therefore, for a given property name, all correlated objects with that name are analyzed together.

The differences between this method of tracking correlated calls and our analysis are the following.

- *Type of target data-flow analysis* whose precision is to be improved. Our analysis improves the precision of IFDS data-flow analyses, whereas the JavaScript analysis improves the precision of pointer analysis.
- *Target language*. Our analysis is for object-oriented languages where polymorphic methods, and not property names (which are known at compile time), cause infeasible paths.
- *Different handling of correlated calls*. Extracting code that contains correlated calls into separate methods would not prevent infeasible paths. Instead, our analysis uses IDE flow functions to detect and eliminate infeasible paths caused by correlated calls.

## 8 Conclusions

We presented a technique to improve the precision of solutions to IFDS problems in the presence of correlated calls. Correlated calls occur when there are multiple polymorphic method invocations on the same receiver. Such method calls cause a data-flow analysis to consider infeasible paths, which makes the data-flow analysis less precise.

Our method of eliminating infeasible paths caused by correlated calls works by transforming an existing IFDS problem into a specialized IDE problem. In this way, we are able to track the classes to which method invocations get dispatched. After solving the specialized IDE problem, we convert its result into an IFDS result that is potentially more precise than the solution to the original IFDS problem. The increase in precision can occur for programs that contain correlated calls. Specifically, if, on a certain data-flow path, there are two polymorphic method invocations on the same receiver that dispatch to incompatible classes, the IDE analysis will consider the path infeasible.

We proved that the correlated-calls analysis is sound and that it improves the precision of IFDS results.

Our Scala implementation of the correlated-calls analysis includes

- an implementation of the IDE analysis, which is based on the WALA static program analysis framework;
- a taint-analysis implementation as an IFDS problem instance;
- a transformer of IFDS problems to equivalent IDE problems, and a second transformer that accounts for correlated calls.

We tested the correlated-calls analysis on our taint analysis implementation by comparing the number of secret values that were leaked when using an IFDS taint analysis and a taint analysis that accounts for correlated calls. We used the DaCapo benchmarks as input programs. Although the benchmarks contained a number of correlated calls, we were not able to improve the precision of the taint analysis, because the correlated calls did not occur on paths of secret information leaks.

We are hopeful that other analyses can benefit from the extra information provided by the correlated-calls analysis, and plan to test this hypothesis in the future.

## References

1. Andrew W. Appel. SSA is functional programming. *SIGPLAN Notices*, (4):17–20, 1998.
2. Stephen M. Blackburn, Robin Garner, Chris Hoffmann, Asjad M. Khan, Kathryn S. McKinley, Rotem Bentzur, Amer Diwan, Daniel Feinberg, Daniel Frampton, Samuel Z. Guyer, Martin Hirzel, Antony L. Hosking, Maria Jump, Han Bok Lee, J. Eliot B. Moss, Aashish Phansalkar, Darko Stefanovic, Thomas VanDrunen, Daniel von Dincklage, and Ben Wiedermann. The DaCapo benchmarks: Java benchmarking development and analysis. In *Proceedings of the 21th Annual ACM*

- SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2006, October 22-26, 2006*, pages 169–190, 2006.
3. Eric Bodden. Inter-procedural data-flow analysis with IFDS/IDE and Soot. In *Proceedings of the ACM SIGPLAN International Workshop on State of the Art in Java Program analysis, SOAP 2012, June 14, 2012*, pages 3–8, 2012.
  4. Eric Bodden, Társis Tolêdo, Márcio Ribeiro, Claus Brabrand, Paulo Borba, and Mira Mezini. SPLIFT - statically analyzing software product lines in minutes instead of years. In *Software Engineering 2014, Fachtagung des GI-Fachbereichs Softwaretechnik, 25. Februar - 28. Februar 2014*, pages 81–82, 2014.
  5. Ron Cytron, Jeanne Ferrante, Barry K. Rosen, Mark N. Wegman, and F. Kenneth Zadeck. Efficiently computing static single assignment form and the control dependence graph. *ACM Trans. Program. Lang. Syst.*, (4):451–490, 1991.
  6. Stephen Fink and Julian Dolby. WALA — the TJ Watson libraries for analysis. <http://wala.sourceforge.net>, 2012.
  7. Laurent Hubert and David Pichardie. Soundly handling static fields: Issues, semantics and analysis. *Electr. Notes Theor. Comput. Sci.*, (5):15–30, 2009.
  8. Jens Knoop and Bernhard Steffen. The interprocedural coincidence theorem. In *Compiler Construction, 4th International Conference on Compiler Construction, CC'92, October 5-7, 1992, Proceedings*, pages 125–140, 1992.
  9. Jens Knoop, Bernhard Steffen, and Jürgen Vollmer. Parallelism for free: Efficient and optimal bitvector analyses for parallel programs. *ACM Trans. Program. Lang. Syst.*, (3):268–299, 1996.
  10. Jörg Kreiker, Thomas W. Reps, Noam Rinetzky, Mooly Sagiv, Reinhard Wilhelm, and Eran Yahav. Interprocedural shape analysis for effectively cutpoint-free programs. In *Programming Logics — Essays in Memory of Harald Ganzinger*, pages 414–445, 2013.
  11. Tim Lindholm and Frank Yellin. *The Java Virtual Machine Specification*. 1997.
  12. Markus Müller-Olm and Oliver Rüthing. On the complexity of constant propagation. In *Programming Languages and Systems, 10th European Symposium on Programming, ESOP 2001 Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2001, April 2-6, 2001, Proceedings*, pages 190–205, 2001.
  13. Nomair A. Naeem and Ondřej Lhoták. Typestate-like analysis of multiple interacting objects. In *Proceedings of the 23rd Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2008, October 19-23, 2008*, pages 347–366, 2008.
  14. Nomair A. Naeem, Ondřej Lhoták, and Jonathan Rodriguez. Practical extensions to the IFDS algorithm. In *Compiler Construction, 19th International Conference, CC 2010, March 20-28, 2010. Proceedings*, pages 124–144, 2010.
  15. Flemming Nielson, Hanne Riis Nielson, and Chris Hankin. *Principles of program analysis (2. corr. print)*. 2005.
  16. Martin Odersky. Essentials of Scala. In *Langages et Modèles à Objets, LMO 2009, 25-27 mars 2009*, page 2, 2009.
  17. Thomas W. Reps, Susan Horwitz, and Shmuel Sagiv. Precise interprocedural dataflow analysis via graph reachability. In *Conference Record of POPL'95: 22nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, January 23-25, 1995*, pages 49–61, 1995.
  18. Henry Gordon Rice. Classes of recursively enumerable sets and their decision problems. *Transactions of the American Mathematical Society*, pages 358–366, 1953.

19. Jonathan David Rodriguez. A concurrent IFDS dataflow analysis algorithm using actors. Master's thesis, 2010.
20. Shmuel Sagiv, Thomas W. Reps, and Susan Horwitz. Precise interprocedural dataflow analysis with applications to constant propagation. In *TAPSOFT'95: Theory and Practice of Software Development, 6th International Joint Conference CAAP/FASE, May 22-26, 1995, Proceedings*, pages 651–665, 1995.
21. Micha Sharir and Amir Pnueli. Two approaches to interprocedural data flow analysis. *Program flow analysis: Theory and applications*, pages 189–234, 1981.
22. Manu Sridharan, Julian Dolby, Satish Chandra, Max Schäfer, and Frank Tip. Correlation tracking for points-to analysis of JavaScript. In *ECOOP 2012 - Object-Oriented Programming - 26th European Conference, June 11-16, 2012. Proceedings*, pages 435–458, 2012.
23. Frank Tip. Infeasible paths in object-oriented programs. *Science of Computer Programming*. To appear, 2014.
24. Omer Tripp, Marco Pistoia, Stephen J. Fink, Manu Sridharan, and Omri Weisman. TAJ: effective taint analysis of web applications. In *Proceedings of the 2009 ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2009, June 15-21, 2009*, pages 87–97, 2009.
25. Raja Vallée-Rai, Phong Co, Etienne Gagnon, Laurie J. Hendren, Patrick Lam, and Vijay Sundaresan. Soot — a Java bytecode optimization framework. In *Proceedings of the 1999 conference of the Centre for Advanced Studies on Collaborative Research, November 8-11, 1999*, page 13, 1999.

## Appendix

In this appendix we present the proofs to the Lemmas introduced in Section 5.

### Soundness and Precision

We start by proving the Lemmas of Soundness and Precision of the correlated-calls analysis.

**Proof of Lemma 1.** The transformation  $\mathcal{U}^\subseteq$  is the same as  $\mathcal{U}^\equiv$ , except that it can remove data-flow facts from the result:

$$\begin{aligned}
 \mathcal{U}^\subseteq(\mathcal{R}(\mathcal{T}_R^\subseteq(P)))(n) &= \{(n', D_n'^\subseteq(\mathcal{R}(\mathcal{T}_R^\subseteq(P)))) \mid n \in N^*\}(n) \\
 &= D_n^\subseteq(\mathcal{R}(\mathcal{T}_R^\subseteq(P))) \\
 &\subseteq \text{MVP}_F(n) \\
 &= \mathcal{R}_{\text{IFDS}}(P)(n). \quad \square
 \end{aligned}$$

To prove the Soundness Lemma, we first introduce Lemmas 7 and 8.

We will denote the top element in the environment lattice as

$$\Omega = \lambda d. \top_\subseteq. \quad (47)$$

For the purpose of the proofs, we will rewrite Equation (23) that defines an edge function as follows:

$$\text{EdgeFn}_S^\subseteq = \lambda e. \begin{cases} \text{id} & \text{if } d_1 = d_2 = \Lambda, \\ \lambda m. \varepsilon(e)(\delta(m)) & \text{otherwise,} \end{cases} \quad (48)$$

where  $S \subseteq R$ ,  $d_1$  and  $d_2$  are the source and target facts, and for a map  $m \in L_U^\subseteq$ ,  $\delta(m)$  is either  $m$  or  $\perp_\subseteq$ :

$$\delta(m) = \begin{cases} \perp_\subseteq & \text{if } d_1 = \Lambda \\ m & \text{otherwise.} \end{cases} \quad (49)$$

Additionally, for a path  $p = [\text{start}_{\text{main}}, \dots]$  and a fact  $d \in D$ , we will denote the lattice element that is mapped to  $d$  according to the flow functions of path  $p$  as follows:

$$\xi(p, d) = M_{\text{Env}}(p)(\Omega)(d). \quad (50)$$

The following Lemma shows that the lattice elements (receiver-to-types maps) of a correlated-calls IDE analysis correctly overapproximate the possible types of a receiver in a program execution.

**Lemma 7.** *Let  $p = [\text{start}_{\text{main}}, \dots, n]$  be some concrete execution trace of the program, and let  $r \in R$  be a receiver. If after the execution trace  $p$ , at node  $n$ ,  $r$  points to an object of runtime type  $t$ , and  $d \in D$  is a fact such that  $d \in M_F(p)(\emptyset)$ , then*

$$t \in \xi(p, d)(r). \quad (51)$$

*Proof.* By induction on the length of the trace.

*Basis:*  $p = [\text{start}_{\text{main}}]$ . Then there is no instruction at which a receiver  $r$  could be instantiated, and the Lemma is trivially true.

*Induction hypothesis:* Let  $p = [\text{start}_{\text{main}}, \dots, n_{k-1}]$ , and let  $\tau$  be the set of types to which  $\xi(p, d_{k-1})$  maps  $r$ :

$$\tau = \xi(p, d_{k-1})(r). \quad (52)$$

Assume that for a concrete execution path  $p = [\text{start}_{\text{main}}, \dots, n_{k-1}]$ , at node  $(n_{k-1}, d_{k-1})$ , the Lemma holds, i.e.  $t \in \tau$ .

*Induction step:* Let  $p' = [\text{start}_{\text{main}}, \dots, n_{k-1}, n_k]$  and  $t' \in T$  be the type to which  $r$  is mapped at  $n_k$ .

For each  $i$ , let  $e_i$  be the edge  $((n_{i-1}, d_{i-1}), (n_i, d_i))$ . Note that

$$e_1 = ((\text{start}_{\text{main}}, \Lambda), (n_1, d_1)).$$

Observe that

$$\begin{aligned} \xi(p', d) &= M_{\text{Env}}(p')(\Omega)(d) \\ &= (M_{\text{Env}}(e_k) \circ M_{\text{Env}}(e_{k-1}) \circ \dots \circ M_{\text{Env}}(e_1))(\Omega)(d) \\ &= M_{\text{Env}}(e_k)(M_{\text{Env}}(p)(\Omega))(d). \end{aligned}$$



According to (18),

$$\begin{aligned}
& M_{\text{Env}}(e_k) (M_{\text{Env}}(p)(\Omega)) (d)(r) \\
&= \left( \text{EdgeFn}_R^{\subseteq}((n_{k-1}, \Lambda), (n_k, d))(\top_{\subseteq}) \right) \sqcap \\
&\quad \left( \bigcap_{d' \in D} \text{EdgeFn}_R^{\subseteq}((n_{k-1}, d'), (n_k, d))(M_{\text{Env}}(p)(\Omega)(d')) \right) (r) \\
&\supseteq \bigcap_{d' \in D} \text{EdgeFn}_R^{\subseteq}((n_{k-1}, d'), (n_k, d))(M_{\text{Env}}(p)(\Omega)(d')) (r) \\
&\supseteq \text{EdgeFn}_R^{\subseteq}((n_{k-1}, d_{k-1}), (n_k, d))(\xi(p, d_{k-1}))(r).
\end{aligned}$$

Therefore,

$$\text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r) \subseteq \xi(p', d)(r). \quad (53)$$

We will now show that

$$t' \in \text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r),$$

which, due to (53), means that the Lemma holds.

According to (48), there are two cases in which  $\text{EdgeFn}_R^{\subseteq}(e_k)$  could fail.

If  $d_{k-1} = d_k = \Lambda$ , then  $d_k \notin M_F(p)(\emptyset)$ , since it does not belong to the set  $D$ , and the Lemma trivially holds.

Otherwise,

$$\text{EdgeFn}_R^{\subseteq}(e_k) = \lambda m. \varepsilon(e_k)(\delta(m)).$$

It follows that

$$\begin{aligned}
\text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r) &= (\lambda m. \varepsilon(e_k)(\delta(m)))(\xi(p, d_{k-1}))(r) \\
&= \varepsilon(e_k)(\delta(\xi(p, d_{k-1}))) (r).
\end{aligned} \quad (54)$$

Let us denote the lattice element  $\delta(\xi(p, d_{k-1}))$  with  $\Delta$ :

$$\Delta = \delta(\xi(p, d_{k-1})).$$

Note that since  $\Delta$ , according to (49), can be either  $\perp_{\subseteq}$  or  $\xi(p, d_{k-1})$ , it always maps  $r$  to a set containing  $t$ :

$$t \in \Delta(r). \quad (55)$$

Note also that unless the instruction at  $n_{k-1}$  contains an assignment for  $r$ ,  $r$  is mapped to the same object of type  $t$  as at node  $n_{k-1}$ , and  $t = t'$ . Therefore, for the non-assignment instructions, it is sufficient to prove that  $t \in \Delta(r)$ .

Depending on the instructions at the nodes  $n_{k-1}$  and  $n_k$ , there are four cases:

1. The instruction at  $n_{k-1}$  is an assignment for a receiver  $r' \in R$ . Since  $\varepsilon_R(e_k) = \lambda m. m[r' \rightarrow \perp_T]$ ,

$$\begin{aligned}
\text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r) &= (\lambda m. m[r' \rightarrow \perp_T])(\Delta)(r) \\
&= \Delta[r' \rightarrow \perp_T](r).
\end{aligned}$$

In the resulting map,  $r'$  is mapped to  $\perp_T$ . Then

- (a) if  $r = r'$ , then  $\text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r) = \perp_T$ , which contains  $t'$ .
  - (b) If  $r \neq r'$ , then  $r$  has not been reassigned a value, and still maps to the same object of type  $t$ . The receiver  $r$  is mapped to  $\Delta(r)$ , which, according to (55), contains  $t$ . Since  $t = t'$ ,  $\Delta(r)$  contains  $t'$ .
2.  $e_k$  is a call-start edge with signature  $s_{\mathcal{F}}$ , and  $f \in \mathcal{F}$  is the called procedure. Then

$$\begin{aligned} \text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r) &= (\lambda m. m[r' \rightarrow m(r') \cap \tau(s_{\mathcal{F}}, f)])(\Delta)(r) \\ &= \Delta[r' \rightarrow \Delta(r') \cap \tau(s_{\mathcal{F}}, f)], \end{aligned}$$

where  $r'$  is the receiver of the call.

- If  $r' = r$ , then  $\Delta(r') = \Delta(r)$  which contains  $t$ . Since  $t \in \tau(s_{\mathcal{F}}, f)$ , it follows that  $t \in \Delta(r) \cap \tau(s_{\mathcal{F}}, f)$ , and  $t \in \text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r)$ .
  - If  $r' \neq r$ , see (1b).
3.  $e_k$  is an end-return edge,  $r_1, \dots, r_k \in R$  are the local variables in the callee method,  $r'$  is the receiver of the call site corresponding to the return node  $n_k$ , and  $f \in \mathcal{F}$  is the called method with signature  $s_{\mathcal{F}}$ . Then

$$\varepsilon_R(e_k) = \lambda m. m[r' \rightarrow m(r') \cap \tau(s_{\mathcal{F}}, f)][r_1 \rightarrow \perp_T] \dots [r_k \rightarrow \perp_T].$$

If  $r \in \{r_1, \dots, r_k\}$ , see Case 1. Otherwise, the case is analogous to Case 2.

4. The node contains any other instruction. Then

$$\text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r) = \text{id}(\Delta)(r) = \Delta(r),$$

which contains  $t$  according to (55). □

We will now show that on a node of a concrete execution path, the correlated-calls analysis does not map receivers to  $\top_T$ . In other words, the analysis never considers nodes of a concrete execution path unreachable.

**Lemma 8.** *Let  $p = [\text{start}_{\text{main}}, \dots, n]$  be a concrete execution path,  $r \in R$  a receiver, and  $d \in D$  a data-flow fact. Then if  $d \in M_F(p)(\emptyset)$ ,*

$$\xi(p, d)(r) \neq \top_T. \quad (56)$$

*Proof.* By induction on the length of the execution trace.

*Basis:* Let  $p = [\text{start}_{\text{main}}]$ . Since the only realizable path corresponding to  $p$  is  $[(\text{start}_{\text{main}}, \Lambda)]$ , there is no fact  $d \in D$  such that  $d \in M_F(p)(\emptyset)$ , and the claim follows immediately.

*Induction hypothesis:* Let  $p = [\text{start}_{\text{main}}, \dots, n_{k-1}]$ . Let  $\tau$  be the set of types to which  $r$  is mapped by  $\xi(p, d_{k-1})$ :

$$\tau = \xi(p, d_{k-1})(r). \quad (57)$$

Assume the Lemma holds for that for a concrete execution path

$$p = [\text{start}_{\text{main}}, n_1, \dots, n_{k-1}],$$

i.e.  $\tau \neq \top_T$  for an arbitrary  $r \in R$  and  $d_{k-1} \in D$ .

*Induction step:* Let  $p' = [\text{start}_{\text{main}}, n_1, \dots, n_{k-1}, n_k]$  be a concrete execution path.

Let  $e_k = ((n_{k-1}, d_{k-1}), (n_k, d))$ . As shown in (53),

$$\xi(p', d)(r) \supseteq \text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r).$$

From Definition 5, we can see that unless  $e_k$  is a call-start edge or an end-return edge, the result follows from the induction hypothesis. More formally, if  $e_k$  is not a call-start or end-return edge, then for all  $m \in L_R^{\subseteq}$ ,

$$\text{EdgeFn}_R^{\subseteq}(e_k)(m) \subseteq m.$$

The edge function corresponding to the call-start and end-return edges is the only place in which the set of types that a receiver maps to can be reduced.

Assume that  $e_k$  is a end-return edge with a call on the receiver  $r' \in R$  with a signature  $s_{\mathcal{F}}$  to a function  $f \in \mathcal{F}$ .

$$\begin{aligned} \text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r) \\ &= (\lambda m. m[r' \rightarrow m(r) \cap \tau(s_{\mathcal{F}}, f)][r_1 \rightarrow \perp_T] \dots [r_l \rightarrow \perp_T])(\xi(p, d_{k-1}))(r) \\ &= (\xi(p, d_{k-1})[r' \rightarrow \tau \cap \tau(s_{\mathcal{F}}, f)][r_1 \rightarrow \perp_T] \dots [r_l \rightarrow \perp_T])(r), \end{aligned}$$

where  $r_1, \dots, r_l \in R$  are the local variables in the called method.

If  $r \in \{r_1, \dots, r_l\}$ , then  $\text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r) = \perp_T \ni t^{15}$ .

Otherwise, if  $r = r'$ , then  $\text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r) = \tau \cap \tau(s_{\mathcal{F}}, f)$ .

According to Lemma 7 and by the induction hypothesis, the runtime type  $t$  of  $r$  must be contained in  $\xi(p, d_{k-1})(r) = \tau$ . At the same time, by definition,  $t$  is part of  $\tau(s_{\mathcal{F}}, f)$ . Therefore,  $t \in \tau \cap \tau(s_{\mathcal{F}}, f) \subseteq \text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r)$ , which means that  $\text{EdgeFn}_R^{\subseteq}(e_k)(\xi(p, d_{k-1}))(r) \neq \top_T$ .

The same reasoning applies to the case where  $e_k$  is a call-start edge.  $\square$

Finally, we can prove the Soundness Lemma.

**Proof of Lemma 2.** Let  $\rho = \mathcal{R}(\mathcal{T}_R^{\subseteq}(P))$ . Then

$$\begin{aligned} \mathcal{U}^{\subseteq}(\rho)(n) &= D_n^{\subseteq}(\rho) \\ &= \{d' \mid d' \in \text{MVP}_F(n) \wedge \forall r \in R: \rho(n, d')(r) \neq \top_T\}. \end{aligned}$$

Since  $\text{MVP}_F(n) = \bigcap_{q \in \text{VP}(n)} M_F(q)(\emptyset)$ , and  $p \in \text{VP}(n)$ , it follows that

$$\begin{aligned} d &\in M_F(p)(\emptyset) \\ &\subseteq \text{MVP}_F(n). \end{aligned}$$

<sup>15</sup> In the case of a recursive call, it is possible that both  $r \in \{r_1, \dots, r_l\}$  and  $r = r'$ . In that case, the set to which  $r$  will be mapped would be still “overwritten” by  $\perp_T$ .

At the same time, for all receivers  $r \in R$ ,

$$\begin{aligned}\rho(n, d)(r) &= \left( \prod_{q \in \text{VP}(n)} \xi(q, d) \right)(r) \\ &= \prod_{q \in \text{VP}(n)} \xi(q, d)(r).\end{aligned}$$

According to Lemma 8,  $\xi(p, d)(r) \neq \top_T$ . Since  $p \in \text{VP}(n)$ ,

$$\xi(p, d)(r) \subseteq \prod_{q \in \text{VP}(n)} \xi(q, d)(r).$$

From  $\prod_{q \in \text{VP}(n)} \xi(q, d)(r) = \rho(n, d)(r)$  it follows that  $\xi(p, d)(r) \subseteq \rho(n, d)(r)$ . Therefore,  $\rho(n, d)(r) \neq \top_T$ , and  $d \in D_n^\subseteq(\rho) = \mathcal{U}^\subseteq(\rho)(n)$ .  $\square$

### Correlated Call Receivers

This appendix contains the proof for Lemma ?? which shows that in a correlated-calls analysis, it is enough to consider only correlated-call receivers  $R^\subseteq$ .

The following Lemma shows that the types to which a given receiver is mapped in the result of the algorithm is not affected by other receivers and the types to which they are mapped.

**Lemma 9.** *Let  $P$  be an IFDS problem. Let  $N^*$  be the supergraph for  $P$ ,  $D$  the set of data-flow facts,  $n \in N^*$  a node, and  $p = [\text{start}_{\text{main}}, \dots, n]$  a path in the supergraph. Let  $d \in D \cup \{\Lambda\}$ . Then for any realizable path  $p' \in \text{RP}(p, d)$ , set  $S \subseteq R$ , and receiver  $r \in S$ ,*

$$\text{EdgeFn}_S^\subseteq(p')(\top_\subseteq)(r) = \text{EdgeFn}_{\{r\}}^\subseteq(p')(\top_\subseteq)(r). \quad (58)$$

*Proof.* By induction on the length of  $p$ .

*Basis:*  $p' = [(\text{start}_{\text{main}}, \Lambda)]$ . Then  $\text{EdgeFn}_S^\subseteq(p') = \text{id} = \text{EdgeFn}_{\{r\}}^\subseteq(p')$ , and the Lemma follows directly.

*Induction hypothesis:* Suppose that for a path  $q = [(\text{start}_{\text{main}}, \Lambda), \dots, (n_{k-1}, d_{k-1})]$ , where  $q \in \text{RP}(n, d)$ , the Lemma holds, i.e. both edge functions map  $r$  to the same set of types  $\tau$ :

$$\begin{aligned}\tau &= \text{EdgeFn}_S^\subseteq(q)(\top_\subseteq)(r) \\ &= \text{EdgeFn}_{\{r\}}^\subseteq(q)(\top_\subseteq)(r).\end{aligned}$$

*Induction step:* Let  $q' = [(\text{start}_{\text{main}}, \Lambda), \dots, (n_{k-1}, d_{k-1}), (n_k, d_k)]$  and the edge  $e_k = ((n_{k-1}, d_{k-1}), (n_k, d_k))$ .

Observe that for any set  $U \subseteq R$  such that  $r \in U$ ,

$$\text{EdgeFn}_U^\subseteq(q')(\top_\subseteq)(r) = \text{EdgeFn}_U^\subseteq(e_k)(\text{EdgeFn}_U^\subseteq(q)(\top_\subseteq))(r). \quad (59)$$

We can see from (48) that there are two cases.

If  $d_{k-1} = d_k = \Lambda$ ,  $\text{EdgeFn}_S^{\subseteq}(e_k) = \text{id} = \text{EdgeFn}_{\{r\}}^{\subseteq}(e_k)$ , and, due to (59),

$$\begin{aligned} \text{EdgeFn}_S^{\subseteq}(q')(\top_{\subseteq})(r) &= \tau \\ &= \text{EdgeFn}_{\{r\}}^{\subseteq}(q')(\top_{\subseteq})(r). \end{aligned}$$

Otherwise, there are four sub-cases.

1.  $e_k$  is a call-start edge,  $r'.c()$  is the call site at  $n_{k-1}$  with signature  $s_{\mathcal{F}}$ ,  $f \in \mathcal{F}$  is the called procedure, and  $r' \in U$ . Then

$$\text{EdgeFn}_U^{\subseteq}(e_k) = \lambda m. \delta(m)[r' \rightarrow \delta(m)(r) \cap \tau(s_{\mathcal{F}}, f)].$$

There are two sub-cases.

- (a) If  $r = r'$ , then, according to (59), the resulting set of types

$$\text{EdgeFn}_U^{\subseteq}(q')(\top_{\subseteq})(r) = \delta(\text{EdgeFn}_U^{\subseteq}(q)(\top_{\subseteq}))(r) \cap \tau(s_{\mathcal{F}}, f).$$

If  $d_{k-1} = \Lambda$ , then  $\delta(\text{EdgeFn}_U^{\subseteq}(q)(\top_{\subseteq}))(r) = \perp_{\subseteq}(r) = \perp_T$ . If  $d_{k-1} \neq \Lambda$ , then  $\delta(\text{EdgeFn}_U^{\subseteq}(q)(\top_{\subseteq}))(r) = \text{EdgeFn}_U^{\subseteq}(q)(\top_{\subseteq})(r) = \tau$ . The set  $\tau(s_{\mathcal{F}}, f)$  is the same for either case.

Therefore, the value of  $\text{EdgeFn}_U^{\subseteq}(q')(\top_{\subseteq})(r)$  has the same result regardless of  $U$ , which means that  $\text{EdgeFn}_S^{\subseteq}(q')(\top_{\subseteq})(r) = \text{EdgeFn}_{\{r\}}^{\subseteq}(q')(\top_{\subseteq})(r)$ , and the Lemma holds.

- (b) If  $r \neq r'$ , then

$$\text{EdgeFn}_U^{\subseteq}(q')(\top_{\subseteq})(r) = \delta(\text{EdgeFn}_U^{\subseteq}(q)(\top_{\subseteq}))(r), \quad (60)$$

which, as we have seen in Case (1a), does not depend on  $U$ , and the Lemma holds.

2.  $e_k$  is an end-return edge,  $r_1, \dots, r_l \in U$  are the local variables in the callee method,  $r'.c()$  is the call corresponding to the return node at  $n_k$ ,  $f \in \mathcal{F}$  is the called method with signature  $s_{\mathcal{F}}$ , and  $r' \in U$ . Then

$$\text{EdgeFn}_U^{\subseteq}(e_k) = \lambda m. \delta(m)[r' \rightarrow \delta(m)(r) \cap \tau(s_{\mathcal{F}}, f)][r_1 \rightarrow \perp_T] \dots [r_l \rightarrow \perp_T].$$

There are three sub-cases.

- (a) If  $r \in \{r_1, \dots, r_l\}$ , then regardless of the value of  $U$ ,

$$\text{EdgeFn}_U^{\subseteq}(q')(\top_{\subseteq})(r) = \perp_T,$$

and the Lemma holds.

- (b) Otherwise, if  $r = r'$ , the case is analogous to Case (1a).
- (c) If  $r \notin \{r', r_1, \dots, r_l\}$ , then see Case (1b).

3.  $n_{k-1}$  contains an assignment for  $r' \in U$ . Then

$$\text{EdgeFn}_U^{\subseteq}(e_k) = \lambda m. \delta(m)[r' \rightarrow \perp_T].$$

If  $r = r'$ , see Case (2a). If  $r \neq r'$ , see Case (1b).

4. Otherwise,

$$\text{EdgeFn}_U^{\subseteq}(e_k) = \lambda m. \delta(m),$$

and the case is analogous to Case (1b).  $\square$

The following Lemma shows that the correlated-calls analysis computes the results for each receiver independently, or separately. To compute the set of types to which a receiver  $r$  is mapped at each exploded-graph node, we can exclude all other receivers in the program from the analysis (recall from (23) that the set of receivers that are considered in the analysis is specified by the set  $S$  in a correlated-calls transformation  $\mathcal{T}_S^{\subseteq}$ ). Therefore, for a given receiver  $r$ , the results of a  $\mathcal{T}_S^{\subseteq}$ - and a  $\mathcal{T}_{\{r\}}^{\subseteq}$ -analysis are the same.

**Lemma 10.** *Let  $P$  be an IFDS problem. Let  $N^*$  be the supergraph for  $P$ ,  $D$  the set of data-flow facts, and  $S \subseteq R$  a set of receivers. Then for any  $n \in N^*$ ,  $d \in D$ , and receiver  $r \in S$ ,*

$$\mathcal{R}(\mathcal{T}_S^{\subseteq}(P))(n, d)(r) = \mathcal{R}(\mathcal{T}_{\{r\}}^{\subseteq}(P))(n, d)(r). \quad (61)$$

*Proof.* According to (??), (9), and (19),

$$\begin{aligned} \mathcal{R}(\mathcal{T}_S^{\subseteq}(P))(n, d)(r) &= \text{MVP}_{\text{Env}}(n, d)(r) \\ &= \left( \bigcap_{q \in \text{VP}(n)} M_{\text{Env}}(q)(\Omega)(d) \right)(r) \\ &= \left( \bigcap_{q \in \text{VP}(n)} \bigcap_{q' \in \text{RP}(q, d)} \text{EdgeFn}_S^{\subseteq}(q')(\top_{\subseteq}) \right)(r) \\ &= \bigcup_{q \in \text{VP}(n)} \bigcup_{q' \in \text{RP}(q, d)} \text{EdgeFn}_S^{\subseteq}(q')(\top_{\subseteq})(r). \end{aligned} \quad (62)$$

Then from Lemma 9,

$$\begin{aligned} \mathcal{R}(\mathcal{T}_S^{\subseteq}(P))(n, d)(r) &= \bigcup_{q \in \text{VP}(n)} \bigcup_{q' \in \text{RP}(q, d)} \text{EdgeFn}_{\{r\}}^{\subseteq}(q')(\top_{\subseteq})(r) \\ &= \mathcal{R}(\mathcal{T}_{\{r\}}^{\subseteq}(P))(n, d)(r). \end{aligned} \quad \square$$

The next lemma shows that the set of types to which a receiver is mapped in a correlated-calls lattice element can be represented as an intersection of static-type function applications  $\tau(s_{\mathcal{F}_i}, f_i)$ .

**Lemma 11.** *For an IFDS problem  $P$ , a node  $n \in N^*$ , and fact  $d \in D$ , let  $p \in \text{RP}(n, d)$  be a realizable path and  $r \in R$  a receiver. Then there exists a non-negative number  $\gamma$  of calls on the receiver  $r$  with signatures  $s_{\mathcal{F}_\gamma}$  to the functions  $f_\gamma \in \mathcal{F}_\gamma$ , for which*

$$\text{EdgeFn}_{\{r\}}^{\subseteq}(p)(\top_{\subseteq})(r) = \bigcap_{\gamma \geq 0} \tau(s_{\mathcal{F}_\gamma}, f_\gamma).$$

*Proof.* Let  $p$  have the following form<sup>16</sup>:

$$p = [(\text{start}_{\text{main}}, \Lambda), (n_1, \Lambda), \dots, (n_k, \Lambda), (n_{k+1}, d_{k+1}), \dots, (n_{k+l}, d_{k+l})],$$

where  $l \geq 1$  and the facts for all nodes up to  $n_k$  are equal to  $\Lambda$  and  $d_{k+i} \in D$  for  $0 < i \leq l$ .

As previously, for all  $i$ , we will denote the edge  $(n_i, n_{i+1})$  by  $e_i$ .

From (23) we can infer that

$$\text{EdgeFn}_{\{r\}}^{\subseteq}(p) = \text{EdgeFn}_{\{r\}}^{\subseteq}(e_{k+l}) \circ \dots \circ \text{EdgeFn}_{\{r\}}^{\subseteq}(e_{k+2}) \circ (\lambda m. \beta) \circ \text{id} \circ \dots \circ \text{id},$$

where

$$\beta = \begin{cases} \perp_{\subseteq}[r \rightarrow \tau(s_{\mathcal{F}}, f)] & \text{if } (n_k, n_{k+1}) \text{ is a call-start or end-return edge, and} \\ & \text{the call site } r.c() \text{ with signature } s_{\mathcal{F}} \text{ to the function} \\ & f \in \mathcal{F} \text{ corresponds to the call-start or end-return edge,} \\ \perp_{\subseteq} & \text{otherwise}^{17}. \end{cases}$$

Therefore,

$$\begin{aligned} \text{EdgeFn}_{\{r\}}^{\subseteq}(p)(\top_{\subseteq}) &= \left( \text{EdgeFn}_{\{r\}}^{\subseteq}(e_{k+l}) \circ \dots \circ \text{EdgeFn}_{\{r\}}^{\subseteq}(e_{k+2}) \right) ((\lambda m. \beta)(\top_{\subseteq})) \\ &= \left( \text{EdgeFn}_{\{r\}}^{\subseteq}(e_{k+l}) \circ \dots \circ \text{EdgeFn}_{\{r\}}^{\subseteq}(e_{k+2}) \circ \text{id} \right) (\beta). \end{aligned} \quad (63)$$

We can now prove the lemma by induction on  $l$ .

*Basis:* If  $l = 1$ , then  $\text{EdgeFn}_{\{r\}}^{\subseteq}(p)(\top_{\subseteq}) = \text{id}(\beta) = \beta$ . There are two cases.

If  $\beta = \perp_{\subseteq}$ , then

$$\begin{aligned} \text{EdgeFn}_{\{r\}}^{\subseteq}(p)(\top_{\subseteq})(r) &= \beta(r) \\ &= \perp_T, \end{aligned}$$

and  $\gamma = 0$ .

If  $\beta = \perp_{\subseteq}[r \rightarrow \tau(s_{\mathcal{F}}, f)]$ , then

$$\text{EdgeFn}_{\{r\}}^{\subseteq}(p)(\top_{\subseteq})(r) = \tau(s_{\mathcal{F}}, f),$$

and  $\gamma = 1$ .

*Induction hypothesis:* Assume that for a path  $p = [(\text{start}_{\text{main}}, \Lambda), \dots, (n_{k+l}, d_{k+l})]$ , the Lemma holds for  $\gamma = N$ , where  $N \geq 0$ .

<sup>16</sup> It can be shown from the definition of a pointwise representation in Sagiv et al. [20] that in a realizable path, there is never an edge from a fact of the set  $D$  to a  $\Lambda$  fact. Therefore, we can represent  $p$  as a sequence of nodes that has a prefix of  $\Lambda$ -fact nodes, after which all nodes are non- $\Lambda$  facts.

<sup>17</sup> Since  $d_k = \Lambda$  and  $d_{k+1} \neq \Lambda$ , the micro function for the edge  $e_{k+1}$  is equal to  $\lambda m. \varepsilon_{\{r\}}(e_{k+1})(\perp_{\subseteq})$ . From the definition of  $\varepsilon_S$  (22) we can see that the only case where  $\varepsilon_{\{r\}}(e_{k+1})(m)$  would not be equal to  $\perp_{\subseteq}$  is when  $e_{k+1}$  is call-start or end-return edge.

*Induction step:* Let  $p' = [(\text{start}_{\text{main}}, \Lambda), \dots, (n_{k+l}, d_{k+l}), (n_{k+l+1}, d_{k+l+1})]$ . Recall that

$$\text{EdgeFn}_{\{r\}}^{\subseteq}(p')(\top_{\subseteq})(r) = \text{EdgeFn}_{\{r\}}^{\subseteq}(e_{k+l+1}) \left( \text{EdgeFn}_{\{r\}}^{\subseteq}(p)(\top_{\subseteq}) \right) (r).$$

From (22) we can see that unless  $e_{k+l+1}$  is a call-start or end-return edge corresponding to a call on the receiver  $r$ , then  $\text{EdgeFn}_{\{r\}}^{\subseteq}(e_{k+l+1})(r)$  must be equal to either  $\perp_T$  or  $m(r)$ , where  $m = \text{EdgeFn}_{\{r\}}^{\subseteq}(p)(\top_{\subseteq})$ .

If  $\text{EdgeFn}_{\{r\}}^{\subseteq}(e_{k+l+1})(r) = \perp_T$ , then the Lemma holds for  $\gamma = 0$ . Otherwise,

$$\begin{aligned} \text{EdgeFn}_{\{r\}}^{\subseteq}(e_{k+l+1})(\top_{\subseteq})(r) &= \text{EdgeFn}_{\{r\}}^{\subseteq}(p)(\top_{\subseteq})(r) \\ &= \bigcap_N \tau(s_{\mathcal{F}_N}, f_N), \end{aligned}$$

and therefore  $\gamma = N$ .

Suppose that  $e_{k+l+1}$  is a call-start edge with a call on the receiver  $r$  with signature  $s_{\mathcal{G}}$  to a function  $g \in \mathcal{G}$ . Then, according to (22),

$$\text{EdgeFn}_{\{r\}}^{\subseteq}(e_{k+l+1}) = \lambda m. m[r \rightarrow m(r) \cap \tau(s_{\mathcal{G}}, g)].$$

Therefore,

$$\begin{aligned} \text{EdgeFn}_{\{r\}}^{\subseteq}(p')(\top_{\subseteq})(r) &= \lambda m. m[r \rightarrow m(r) \cap \tau(s_{\mathcal{G}}, g)] \left( \text{EdgeFn}_{\{r\}}^{\subseteq}(p)(\top_{\subseteq}) \right) (r) \\ &= \text{EdgeFn}_{\{r\}}^{\subseteq}(p)(\top_{\subseteq})(r) \cap \tau(s_{\mathcal{G}}, g) \\ &= \left( \bigcap_N \tau(s_{\mathcal{F}_N}, f_N) \right) \cap \tau(s_{\mathcal{G}}, g), \end{aligned}$$

and the Lemma holds for  $\gamma = N + 1$ .

The case where  $e_{k+l+1}$  is an end-return edge is analogous to the previous case.  $\square$

We now show that a receiver will be only mapped to  $\top_{\subseteq}$  if it is the receiver of a correlated call.

**Lemma 12.** *For an IFDS problem  $P$ , let  $n \in N^*$  be a node, and  $d \in D$  a dataflow fact such that there exists a realizable path  $p \in \text{RP}(n, d)$ . Let  $T$  be the set of all types in the program. If there exists a receiver  $r \in R$  such that*

$$\text{EdgeFn}_{\{r\}}^{\subseteq}(p)(\top_{\subseteq})(r) = \top_T,$$

*then  $r \in R^{\subseteq}$ .*



*Proof.* According to Lemma 11,

$$\text{EdgeFn}_{\{r\}}^{\subseteq}(p)(\top_{\subseteq})(r) = \bigcap_{\gamma \geq 0} \tau(s_{\mathcal{F}_\gamma}, f_\gamma).$$

Let  $\tau_i = \tau(s_{\mathcal{F}_i}, f_i)$ . For a given  $k$ , let  $r.m_k()$  be the call site corresponding to  $\tau_k$ , and  $T'$  the set of types compatible with the static type of  $r$ . Recall from Section 3.3 that

- $\tau_k \neq \top_T$ ;
- if  $\tau_k = T'$  then the corresponding call site is monomorphic;
- if  $\tau_k \subset T'$  then the call site is polymorphic.

From the conditions of the Lemma,

$$\bigcap_{\gamma \geq 0} \tau_\gamma = \top_T. \quad (64)$$

If all  $\tau_k = T'$ , then  $\bigcap_{\gamma \geq 0} \tau_\gamma$  is also equal to  $T'$ . Since  $T' \neq \top_T$ , this is a contradiction.

If exactly one  $\tau_k \subset T'$  and the rest are equal to  $T'$ , then  $\bigcap_{\gamma \geq 0} \tau_\gamma$  is equal to  $\tau_k$ , which cannot be  $\top_T$  either.

Therefore, there are at least two sets,  $\tau_i$  and  $\tau_j$ , which are strict subsets of  $T'$ . Since both  $\tau_i$  and  $\tau_j$  are non-empty and their intersection equals  $\top_T$ ,  $\tau_i$  and  $\tau_j$  must be disjoint. If  $\tau_i$  and  $\tau_j$  are disjoint, they must correspond to different call sites.

In other words, there are at least two calls on the same receiver for which the static-type function is a strict subset of the set of types compatible with a given receiver  $r$ . It follows that both calls have to be polymorphic. Therefore,  $r \in R^{\subseteq}$ .  $\square$

We will now show that if a receiver ever gets mapped to top, then it is a correlated-calls receiver.

**Lemma 13.** *For an IFDS problem  $P$ , let  $n \in N^*$  be a node, and  $d \in D$  a dataflow fact such there exists a realizable path  $p \in RP(n, d)$ . Then, if there exists a receiver  $r \in R$ , such that*

$$\mathcal{R}(\mathcal{T}_{\{r\}}^{\subseteq}(P))(n, d)(r) = \top_T,$$

*then  $r \in R^{\subseteq}$ .*

*Proof.* As shown in (62),

$$\mathcal{R}(\mathcal{T}_{\{r\}}^{\subseteq}(P))(n, d)(r) = \bigcup_{q \in \text{VP}(n)} \bigcup_{q' \in \text{RP}(q, d)} \text{EdgeFn}_{\{r\}}^{\subseteq}(q')(\top_{\subseteq})(r).$$

Since the latter is equal to  $\top_T$ , it follows that for each realizable path  $p'$  to node  $n$ ,  $\text{EdgeFn}_{\{r\}}^{\subseteq}(p')(\top)(r) = \top_T$ . According to Lemma 13, this is only possible if  $r \in R^{\subseteq}$ .  $\square$

Finally, we present the proof for Lemma ?? which states that a correlated-calls analysis that considers all receivers computes the same result as an analysis that considers only correlated-call receivers.

**Proof of Lemma 3.** From (??) we know that

$$\mathcal{U}^{\mathbb{E}}(\mathcal{R}(\mathcal{T}_R^{\mathbb{E}}(P))) = \{(n, D_n^{\mathbb{E}}(\mathcal{R}(\mathcal{T}_R^{\mathbb{E}}(P)))) \mid n \in N^*\}.$$

According to (??) and Lemma 10, for a given  $n \in N^*$ ,

$$\begin{aligned} D_n^{\mathbb{E}}(\mathcal{R}(\mathcal{T}_R^{\mathbb{E}}(P))) &= \left\{ d \mid d \in \text{MVP}_F(n) \wedge \forall r \in R : \left\{ (r, \mathcal{R}(\mathcal{T}_{\{r\}}^{\mathbb{E}}(P))(n, d)(r)) \mid r \in R \right\} (r) \neq \top_T \right\} \\ &= \left\{ d \mid d \in \text{MVP}_F(n) \wedge \forall \mathbf{r} \in \mathbf{R} : \mathcal{R}(\mathcal{T}_{\{\mathbf{r}\}}^{\mathbb{E}}(P))(n, d)(\mathbf{r}) \neq \top_T \right\}. \end{aligned}$$

Since, according to Lemma 13,  $\mathcal{R}(\mathcal{T}_{\{r\}}^{\mathbb{E}}(P))(n, d)(r)$  can only be equal to  $\top_T$  when  $r \in R^{\mathbb{E}}$ , we can conclude that

$$\begin{aligned} D_n^{\mathbb{E}}(\mathcal{R}(\mathcal{T}_R^{\mathbb{E}}(P))) &= \left\{ d \mid d \in \text{MVP}_F(n) \wedge \forall \mathbf{r} \in \mathbf{R}^{\mathbb{E}} : \mathcal{R}(\mathcal{T}_{\{\mathbf{r}\}}^{\mathbb{E}}(P))(n, d)(\mathbf{r}) \neq \top_T \right\} \\ &= D_n^{\mathbb{E}}(\mathcal{R}(\mathcal{T}_{R^{\mathbb{E}}}^{\mathbb{E}}(P))). \end{aligned}$$

Therefore,

$$\begin{aligned} \mathcal{U}^{\mathbb{E}}(\mathcal{R}(\mathcal{T}_R^{\mathbb{E}}(P))) &= \{(n, D_n^{\mathbb{E}}(\mathcal{R}(\mathcal{T}_{R^{\mathbb{E}}}^{\mathbb{E}}(P)))) \mid n \in N^*\} \\ &= \mathcal{U}^{\mathbb{E}}(\mathcal{R}(\mathcal{T}_{R^{\mathbb{E}}}^{\mathbb{E}}(P))). \end{aligned} \quad \square$$

## Representation of Micro Functions

In this part of the appendix we present the proofs to the Lemmas related to the representation of micro functions with update maps.

**Proof of Lemma 4.** Let  $\text{update}_{f,r}^* = \langle I, U \rangle$ . For any  $\tau \in T$ ,

$$\begin{aligned} \llbracket \mathcal{N}(\text{update}_{f,r}^*) \rrbracket (\tau) &= \llbracket \mathcal{N}(\langle I, U \rangle) \rrbracket (\tau) \\ &= \llbracket \langle I \cup U, U \rangle \rrbracket (\tau) \\ &= \tau \cap (I \cup U) \cup U \\ &= (\tau \cap I) \cup (\tau \cap U) \cup U \\ &= \tau \cap I \cup U \\ &= \llbracket \langle I, U \rangle \rrbracket (\tau) \\ &= \llbracket \text{update}_{f,r}^* \rrbracket (\tau). \end{aligned}$$

Thus,  $\llbracket \mathcal{N}(\text{update}_{f,r}^*) \rrbracket = \llbracket \text{update}_{f,r}^* \rrbracket$ .  $\square$

**Proof of Lemma 5.** Let us show that there always exists a set  $\tau \subseteq T$  such that  $\llbracket \langle I, U \rangle \rrbracket(\tau) \neq \llbracket \langle I', U' \rangle \rrbracket(\tau)$ . There are two cases:

1.  $U \neq U'$ . Then for the empty set  $\tau = \emptyset$ ,

$$\llbracket \langle I, U \rangle \rrbracket(\tau) = \llbracket \langle I, U \rangle \rrbracket(\emptyset) = (\emptyset \cap I) \cup U = U,$$

whereas

$$\llbracket \langle I', U' \rangle \rrbracket(\tau) = \llbracket \langle I', U' \rangle \rrbracket(\emptyset) = (\emptyset \cap I') \cup U' = U'.$$

Hence,  $\llbracket \langle I, U \rangle \rrbracket \neq \llbracket \langle I', U' \rangle \rrbracket$ .

2.  $I \neq I'$ . Then for the set of all types  $\tau = T$ ,

$$\llbracket \langle I, U \rangle \rrbracket(\tau) = \llbracket \langle I, U \rangle \rrbracket(T) = (T \cap I) \cup U = I \cup U.$$

Since  $\langle I, U \rangle$  is normalized,  $U \subseteq I$ , and

$$I \cup U = I.$$

At the same time,

$$\llbracket \langle I', U' \rangle \rrbracket(\tau) = \llbracket \langle I', U' \rangle \rrbracket(T) = (T \cap I') \cup U' = I' \cup U' = I'.$$

Since  $I \neq I'$ , it follows that  $\llbracket \langle I, U \rangle \rrbracket \neq \llbracket \langle I', U' \rangle \rrbracket$ . □

**Proof of Lemma 6.** 1. The identity function is represented as

$$\llbracket \text{id} \rrbracket = \{(r, \langle \perp_T, \top_T \rangle) \mid r \in R^\subseteq\};$$

the top function is represented as

$$\llbracket \lambda m. \top_\subseteq \rrbracket = \{(r, \langle \top_T, \top_T \rangle) \mid r \in R^\subseteq\}.$$

2. Equations (41) and (42) show that the representation of micro functions is closed under composition and meet.
3. To show that our representation for micro functions forms a lattice with finite height, let us first show that  $L_{R^\subseteq}^\subseteq : R^\subseteq \rightarrow 2^T$  forms a lattice. Since  $T$  is a finite set,  $(2^T, \subseteq)$  is a finite-height lattice.  $R^\subseteq$  is a finite set. Hence, the mapping

$$R^\subseteq \mapsto 2^T = \{(r, t) \mid r \in R^\subseteq, t \in 2^T\} = L_{R^\subseteq}^\subseteq$$

also forms a finite-height lattice [15].

Furthermore,  $L_{R^\subseteq}^\subseteq$  is a finite set. Every element of  $L_{R^\subseteq}^\subseteq$  can be applied to  $|R^\subseteq|$  receivers, where each receiver is mapped to a set of types. There are  $|R^\subseteq| \cdot 2^{|T|}$  different possibilities to form those mappings, so

$$|L_{R^\subseteq}^\subseteq| = |R^\subseteq| \cdot 2^{|T|}.$$

Therefore,  $L_{R^\subseteq}^\subseteq \mapsto L_{R^\subseteq}^\subseteq$  also forms a finite-height lattice.

4. All operations can be computed in  $O(R^\subseteq \times T)$  time. Note that the  $R^\subseteq$  and  $T$  sets are an input to the correlated-calls analysis, and the time it takes to compute the meet or composition of micro functions is independent of the representation of the specific operand micro functions.
5. The space bound is  $O(R^\subseteq \times T)$ . □