

# IDEs-Friendly Interprocedural Analyser

Ilya Nozhkin  
Saint Petersburg State University  
St. Petersburg, Russia  
nozhkin.ii@gmail.com

Semyon Grigorev  
Saint Petersburg State University  
St. Petersburg, Russia  
semen.grigorev@jetbrains.com

## ABSTRACT

TODO: ABSTRACT

### ACM Reference Format:

Ilya Nozhkin and Semyon Grigorev. 2019. IDEs-Friendly Interprocedural Analyser. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 3 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Static analyses are important part of modern development tools. They take care of verifying correctness of some program's behaviour freeing a programmer from this duty. By used scope of program, an analyze can be classified as intraprocedural or interprocedural, i.e. as those which make decisions based on only one current procedure or based on the whole program respectively. And interprocedural analyses, in theory, can be more precise due to amount of available information.

```
class A {
    [Source]
    int Source;
}

class B {
    [Filter]
    static int Filter(int d);
}

class C {
    [Sink]
    void Sink(int d);
}

class D {
    void Process(A a) {
        int d = Read(a);
        int f = B.Filter(d);
        Consume(d);
        AnotherConsume(f);
    }

    int Read(A a) {
        return a.Source;
    }

    void Consume(int d) {
        C c = new C();
        c.Sink(d);
    }

    void AnotherConsume(int d) { ... }
}
```

Figure 1: Sample code

For example, let's consider the listing in C# at fig. 1. It is known that method *Sink* is vulnerable to invalid arguments and so, it is marked by appropriate attribute. The method

*Filter* validates its argument and possibly modifies it to ensure that returned data is definitely valid. Methods with such behaviour are marked with attribute *Filter*. The most common example of such method is the one that adds escape characters into the string to avoid affecting internal behaviour of the system by user input. And the field *Source* is known as potentially tainted that is also indicated by attribute. Class *D* extracts data using class *A* as a source, then validates it using class *B* then performs some computations involving class *C*. So, there might appear an issue that leads to usage of data that are not been validated yet as it happens during the invocation of method *Consume*. The problem is to find all such issues, i.e. we want an analysis that finds all possible ways how data from source can reach sink bypassing filters.

This problem is a special case of label-flow analysis, so there are several approaches of solving such problems. One of them is CFL-reachability. (TODO: ADV: performance, DIS-ADV: expressive power, non-obvious structure, CITATIONS) Another is abstract interpretation. (TODO: vice versa) We propose to combine these two approaches to achieve acceptable performance and expressive power. I.e. the program is translated into a graph as it is in CFL-r, but constraints that specify what paths is needed to be accepted are set by pushdown automaton which transition relation simulates the semantics of original program. Let's take a closer look at these two components that define an analysis in conjunction.

## 2 ANALYSIS DEFINITION

### 2.1 Graph extraction

The graph that is explored during analysis is an aggregate of control-flow graphs of each method. The one that corresponds to our example is shown at fig. 2.

Each edge contains an operation that represents a statement in the source code and in the same time the target of the edge indicates where to jump after execution of the operation. For example, there exist three different types of operations: invocations, assignments and returns. Each of them is an image of some source instruction. Invocations are produced from call sites and have the same information as ones in original code. I.e. it contains a reference to the entity which method is called, the name of target, a set of arguments and the information about returned values. Each assignment corresponds to real assignment and has references to the source entity and the target one. And return just indicates the end of a method. It can be not present explicitly in the source code but is still needed to be added to inform analyser about return point. However, the number of instruction types is not fixed and some analysis-specific

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Conference'17, July 2017, Washington, DC, USA

© 2019 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM. . . \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

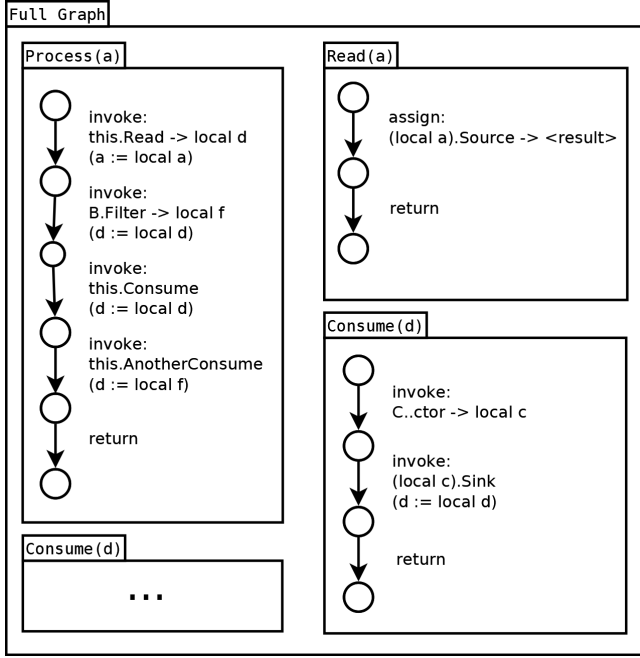


Figure 2: Sample graph

operations can be added if necessary. Nodes have no any data and correspond to positions between instructions in the source code.

So, we have a bunch of graphs each of which represents the content of one method. Next step is to interconnect them to have an opportunity to perform interprocedural jumps during invocations. There are several possible way to do that. First of them is to expand invocations statically, i.e. add a pair of edges for each target of each invocation. One to represent a jump from the call site to the entry point of target and one to emulate return from the final node of the method to the caller. But this approach leads to the need to update all these additional connections if some method is removed or its body is changed. So, we propose to resolve invocations dynamically right during an analysis. It allows us not to modify initial graphs structure at all. Nevertheless, it is still needed to have a mechanism that can perform the resolution, i.e. there must be an entity that can collect all targets of any invocation using references stored there.

To implement such mechanism, called resolver, we offer to accumulate some meta-information about the program besides graphs themselves. The relations in required data is shown at fig. 3.

Firstly, it is important to keep the hierarchy of inheritance to support polymorphic calls and invocations of methods of a basic class. Secondly, it is needed to know which methods are contained in each class to find the method by its name and its location. Thirdly, methods can have local functions and it is necessary to keep their hierarchy too to support, for example, anonymous function invocations, delegates passing and so on. And finally, methods themselves has references

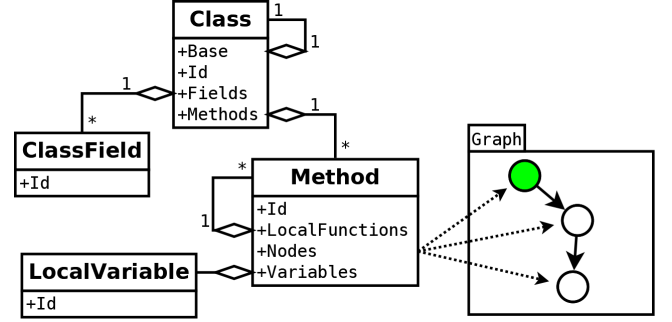


Figure 3: Metadata

to nodes they own which is used to find the entry point and update the graph when the body of method is changed. This structure also contains such data as class fields and local variables of a method which can be referenced by operations. So, the resolver takes the class name and the identifier of a method or a field and walks through the hierarchy tying to find all suitable entities.

## 2.2 PDA construction

Formally, nondeterministic pushdown automaton is a tuple  $(Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$ , where  $Q$ ,  $\Sigma$  and  $\Gamma$  are finite sets of states, input symbols and stack symbols respectively,  $q_0 \in Q$  and  $Z_0 \in \Gamma$  are initial state and stack symbol,  $F \subseteq Q$  is a set of final states and  $\delta : Q \times \Sigma \cup \{\epsilon\} \times \Gamma \rightarrow \mathcal{P}(Q \times \Gamma^*)$  is a transition relation which computes new state and stack by current ones and input symbol. We also add following restriction on transition relation. The resulting stack must differ from the source one by no more than one top symbol. I.e. only one symbol can be pushed or popped during the transition.

Next, we propose to take the set of all edges in the control-flow graph as  $\Sigma$ . So, the transition relation can be understood as a structural operational semantics that defines how configuration is changed during execution of a statement. All other sets can be taken arbitrary.

However, there is one more problem. Semantics of invocation contains the need to make a jump from the current position to the entry point of the target instead of just going to the node that is pointed by the current edge. To support such behaviour we propose to change the codomain of  $\delta$  such that  $\delta : Q \times \Sigma \cup \{\epsilon\} \times \Gamma \rightarrow \mathcal{P}(Q \times \Gamma^* \times N \cup \{\nu\})$ , where  $N$  is the set of nodes of the graph and  $\nu$  is the dummy value that means that there is no need to jump and PDA just goes to the next node.

For example, let's construct the PDA performing that taint tracking analysis described in the introduction. Let  $Q := V \cup \{q_0, q_f\}$ , where  $V$  is set of all local variables of all methods and  $q_0$  and  $q_f$  are dummy initial and final state, so  $F := q_f$ .  $\Gamma := I \cup \{Z_0\}$ , where  $I \subset \Sigma$  is set of all edges containing an invocation and  $Z_0$  is dummy initial stack symbol. And  $\delta$  is

defined by the case analysis (1).

$$\begin{aligned}
\delta(q_0, i@invocation, \gamma) &:= \\
&\quad \{(q_0, i :: \gamma, s_0), \dots, (q_0, i :: \gamma, s_n) : \\
&\quad \quad s_0, \dots, s_n \in R(i)\} \\
\delta(q_0, a@assignment(v_s, v_t), \gamma) &:= \\
&\quad \begin{cases} \{(v_t, \gamma, \nu)\}, & \text{if } source(v_s) \\ \{(q_0, \gamma, \nu)\}, & \text{otherwise} \end{cases} \\
\delta(v, a@assignment(v, v_t), \gamma) &:= (v_t, \gamma, \nu) \\
\delta(v, i@invocation, \gamma) &:= \\
&\quad \bigcup_{j=0}^n \{(v_{j0}, i :: \gamma, s_j), \dots, (v_{jm}, i :: \gamma, s_j) : \\
&\quad \quad v_{jk} \in A(i, j, v), s_j \in R(i)\} \\
\delta(v, r@return, i :: \gamma) &:= \\
&\quad \begin{cases} \{(RV(i), \gamma, T(i))\}, & \text{if } returned(v) \\ \{\}, & \text{otherwise} \end{cases} \\
\delta(q, -, \gamma) &:= \{(q, \gamma, \nu)\}
\end{aligned} \tag{1}$$

Where *source* checks if a variable is a source, *returned* checks if current variable is a return value of some method, *T* returns target node of an edge, *RV* returns the variable where the result of an invocation must be put, *R* is the resolver returning entry points of all possible targets of an invocation and *A* is defined by equation (2).

$$A(i, j, v) := \begin{cases} \{q_f\}, & \text{if } j\text{-th target of invocation } i \\ & \text{is sink and } v \text{ is its argument} \\ \{v_k : v \mapsto v_k\}, & \text{if } j\text{-th target of } i \\ & \text{is not filter} \\ \{\}, & \text{otherwise} \end{cases} \tag{2}$$

Where  $v \mapsto v_k$  means that  $v$  is passed as  $k$ -th parameter and becomes local variable  $v_k$  of the target after passing.

### 3 SOLUTION

The solution is an extensible infrastructure which is responsible for extracting graphs from the source code, aggregating them and their metadata into one database and finding paths in this database accepted by PDAs representing different analyses. Logically, it is divided into two separate entities which are shown on fig. 4.

The first entity, the core of the solution, is a backend implemented as a remote service running in a separate process and interacting with the frontend using a socket-based protocol. Architecturally, it is also divided into two subsystems. First of them is a database which provides the continuous incremental updating of the graph and its metadata, and supports dumping to a disk and further restoration in the beginning of next session. The second is responsible for execution of analyses. It contains an implementation of the resolver, the algorithm of PDAs running and the first extension point making the adding new analyses possible. The set of analyses contained in the backend can be extended by adding a new PDA as just a new implementation of appropriate generic abstract class. Further, it is possible to run this new analysis

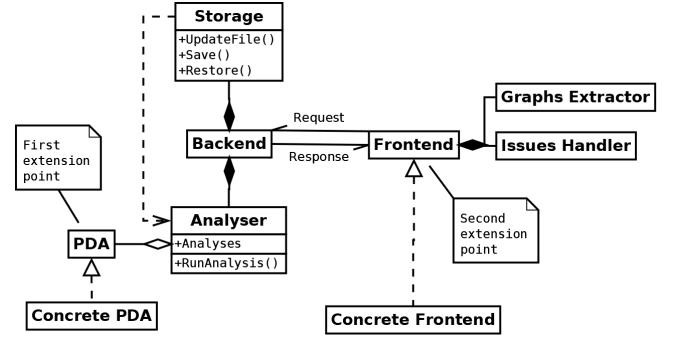


Figure 4: Solution structure

using existing internal algorithm of PDA simulation and get any finite subset of paths in the graph which are accepted by the PDA.

Second main entity is a frontend that is also divided into two subsystems. First of them is a graph extractor which parses source code, extracts graphs and metadata from it and sends collected data to the backend. The second is a results interpreter that receives the set of paths in the graph leading to an error, maps it into the source code and does something with this information. For example, it can highlight pieces of code which participate in the error producing. Also, the frontend is the second extension point because the backend, in general, can interact with any implementation supporting the protocol.

The protocol itself is based on request-response behaviour where the frontend acts as a master and the backend acts as a slave. I.e. the frontend informs the backend if there are some changes in the source code and asks it to update the database according to them. And when there is a need to get the results of some analysis, for example, when the IDE performs the code highlighting, the frontend asks the backend for found issues and waits until it responses.

### 4 EVALUATION

TODO: EVALUATION

### 5 CONCLUSION

TODO: CONCLUSION