

CFL-Reachability Based Framework for Interprocedural Static Code Analysis Development

Ilya Nozhkin

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

Semyon Grigorev

Saint Petersburg State University
St. Petersburg, Russia
JetBrains Research
St. Petersburg, Russia
s.v.grigoriev@spbu.ru
semen.grigorev@jetbrains.com

ABSTRACT

We propose an extensible framework for interprocedural static code analysis implementation. Our solution is based on CFL-reachability: analysis is formulated in terms of context-free constrained reachability in the interprocedural graph. Extensible architecture allows one to implement new analysis and integrate it into the IDE of choice or static code analysis tool. To demonstrate the abilities of our solution, we implement the plugin which provides basic taint analysis and variable flow analysis upon ReSharper infrastructure. We demonstrate its applicability for real-world problems.

CCS CONCEPTS

• **Software and its engineering** → **Automated static analysis; Integrated and visual development environments; Software defect analysis; Software safety**; • **Theory of computation** → **Grammars and context-free languages**;

KEYWORDS

Static code analysis, interprocedural analysis, CFL-reachability, taint analysis, IDE, plugin, context-free languages, PDA

ACM Reference Format:

Ilya Nozhkin and Semyon Grigorev. 2019. CFL-Reachability Based Framework for Interprocedural Static Code Analysis Development. In *Proceedings of Central and Eastern European Software Engineering Conference Russia (CEE-SECR'19)*. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Static code analysis is an important part of modern software development tools. It takes care of verifying the correctness of some program's behavior freeing a programmer from this duty. Static analysis can be classified as intraprocedural or interprocedural depending on which scope of a program it uses. Intraprocedural analysis considers only one procedure whereas interprocedural

inspects a program as a whole. Interprocedural analyses, in theory, are more precise due to the amount of available information.

```
File 1:
1 class A {
2     [Tainted]
3     int Source;
4 }
5
6 class B {
7     [Filter]
8     static int Filter(int d);
9 }
10
11 class C {
12     [Sink]
13     void Sink(int d);
14 }

File 2:
1 class D {
2     void Process(A a) {
3         int d = Read(a);
4         int f = B.Filter(d);
5         Consume(d);
6         AnotherConsume(f);
7     }
8
9     int Read(A a) {
10         return a.Source;
11     }
12
13     void Consume(int d) {
14         C c = new C();
15         c.Sink(d);
16     }
17
18     void AnotherConsume(int d) { ... }
19 }
```

Figure 1: Sample code in C# which represents a case for interprocedural taint analysis

One classical problem which requires interprocedural analysis is taint tracking problem (e.g. described in [2]). Input data can have an inappropriate format or contain an exploit such as SQL injection. Such data is called *tainted*. Tainted data can lead to incorrect behavior in case of incorrect format, or can cause a security issue.

A simple example of C# code which may require taint analysis is presented in fig. 1. We use this example to illustrate our solution. First of all, we introduce a number of entities which play important roles in analysis. The first is *source*: a field that potentially contains the tainted data, for example, the field *Source* of the class *A*. The second is *sink*: a method which is vulnerable to tainted data. In our example, the method *Sink* of the class *C* has this property. The third important type of entities is *filter* (or *sanitizer* in some other definitions). It is a method that checks the correctness of data passed into it. If data is incorrect, *filter* throws an exception or modifies data to ensure the correctness of the result. The method *Filter* of the class *B* in the given snippet is a filter.

We assume that each entity is annotated by a programmer with an appropriate attribute: *[Tainted]* for sources, *[Filter]* for filters and *[Sink]* for sinks. The taint analysis problem is to find all traces such that data can flow into a sink bypassing any filter.

In our example the class *D* extracts data using the class *A* as a source than validates it using class *B*, after that performs some computations involving the class *C*. Invocation of the method *Consume*

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.

CEE-SECR'19, November 14–15(16), 2019, St.Petersburg, Russia

© 2019 Association for Computing Machinery.

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

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

leads to using the data which have not been properly validated. To identify such issues we need to run the taint analysis.

One of the well-known frameworks for interprocedural static code analysis is a CFL-reachability framework proposed by Thomas Reps [8]. This framework is generic: it provides the abstraction which allows one to implement many different types of interprocedural static code analysis, such as pointer analysis [13, 15], taint analysis [5], label flow analysis [6, 7], library summarization [11]. Moreover, CFL-reachability is a long-time studied framework and thus there is a number of specific solutions which demonstrate reasonable performance in practice [12].

It is important for static code analysis to be flow-sensitive that means correct handling call-return pairs: when we have an interprocedural graph, we have no explicit synchronization of call and return edges. As a result, one can find a path which contains incorrect sequences of calls and returns. Also, it is important to provide correct handling of read-write, get-set and other pairs of edges (operations). Note, that CFL-reachability framework allows one to implement analysis which can be flow-sensitive or context-sensitive, but not both [9]. In order to provide both flow- and context-sensitivity one can use more powerful framework [14]. In our examples, we implement flow-sensitive analysis, but one can modify it to be context-sensitive.

The main idea of CFL-reachability is to find paths in a graph that satisfy constraints defined by a context-free grammar. In particular, a path is accepted if the concatenation of labels on its edges forms a word which can be derived in the grammar. However, in practice, there are few drawbacks of such a grammar-based definition. Firstly, grammar-driven parsing is based on exact matching of terminals which forces to generate a very large grammar in the case when edges contain some unique attributes. For example, brackets matching described in [7] or [13, 15] requires to generate as many rules as there are call sites in the source code. In our example, to handle all call-return pairs in paths correctly, the grammar should contain the following productions.

$$\begin{aligned}
 & \dots \\
 & B \rightarrow (2:3 \ B)_{2:3} \ B \\
 & B \rightarrow (2:4 \ B)_{2:4} \ B \\
 & B \rightarrow (2:5 \ B)_{2:5} \ B \\
 & B \rightarrow (2:6 \ B)_{2:6} \ B \\
 & B \rightarrow (2:15 \ B)_{2:15} \ B \\
 & \dots \\
 & B \rightarrow \varepsilon \\
 & \dots
 \end{aligned} \tag{1}$$

Where (file:line) means a call from the specific position and $_{\text{file:line}}$ means a return to the specific position in the code.

Secondly, grammar provides an additional level of abstraction. This level is useful for algorithms formalization, but it is not as flexible as direct manipulation with correspondent push-down automaton.

CFL-reachability requires a graph representing the whole program. It is not feasible to construct such a graph from scratch each time a programmer modifies a program. We need to ensure that the

graph is modified with minial effort. To create a generic framework we also need to make sure the graph representation is generic.

The main goal of our work is to implement a tool solving all mentioned problems and thus allowing to use CFL-reachability in practice. We make the following contributions in the paper.

- We describe the scalable representation of a graph which allows containing as much information about the program as necessary.
- We introduce the approach for the definition of constraints based on pushdown automata instead of grammars which eases formulation of analyses.
- We present the implementation of the proposed approach which allows one to create new types of analysis using introduced abstractions. Also we implement¹ a plugin which uses the solution to provide analysis results to ReSharper, Rider, and InspectCode.
- We evaluate the implementation on both synthetic tests and large open source projects.

2 ANALYSIS DEFINITION

To apply the CFL-reachability framework, the first thing to do is to define the representation of a graph and path restrictions. We assume that a graph is an image of a program and it stores the information about the source code that can be mapped back to locate issues in the sources which are found by the analysis of the graph. Then restrictions on paths define the sequences of operations leading to an issue.

In order to keep the expressive power of the CFL-reachability approach but make the implementation easier, we propose to use pushdown automata instead of grammars which still have equivalent power [4]. Further, we take a closer look at each component and consider the construction of them in application to our example.

2.1 Graph extraction

The graph that is explored during the analysis is the union of control flow graphs each of which is extracted from one distinct method. The graph which corresponds to our example is shown in fig. 2.

Each edge contains an operation that represents a statement in the source code and at the same time, the target of the edge indicates the target to jump after the execution of the operation. We classify operations into a few types and operations which have the same type differs only by unique attributes. However, the number of types is not fixed and some analysis-specific operations can be added if necessary. In our example, we consider three different types of operations: invocations, assignments, and returns. Each of them is an image of some source code instruction. Invocations are produced from call sites and have the same information as ones in the original code. Their notation has the following form:

$$\begin{aligned}
 & \text{invoke:} \\
 & o.m \rightarrow v \\
 & (f_1 := c_1; \dots; f_k := c_k)
 \end{aligned} \tag{2}$$

Where o is an object or a class which method is called, m is the name of a method, v is a variable where the result is stored and f_i, c_i are pairs of the formal and the actual parameter respectively.

¹Source code and executables can be downloaded here: github.com/gsvgit/CoFRA.

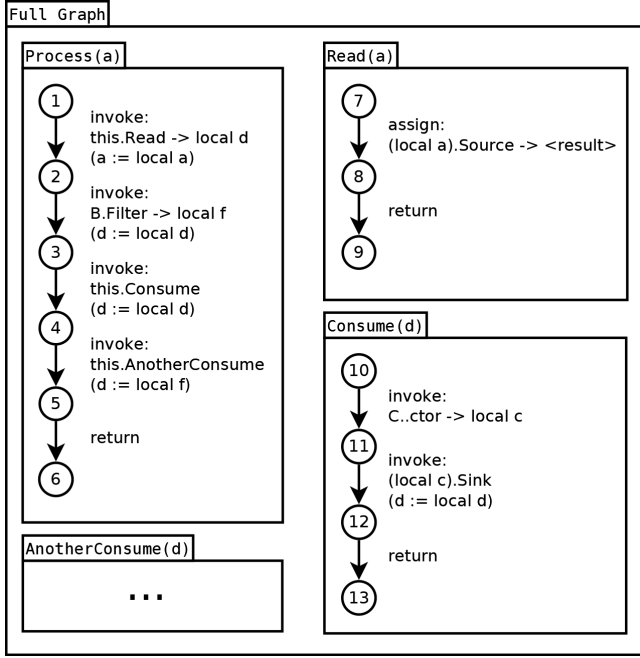


Figure 2: Extracted graph for the code presented in fig. 1

Each assignment corresponds to a real assignment and is written in the following way:

$$\begin{array}{l} \text{assign:} \\ s \rightarrow t \end{array} \quad (3)$$

Where s is a source of data and t is the target variable.

And a return statement 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 the analyzer about the return point.

Nodes have no 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 ways 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 to use the graph which is composed right of graphs of methods and has no additional edges. Also, it does not require to modify any other method when someone is updated. Nevertheless, it is still needed to have a mechanism that can collect all targets of any invocation using references stored there.

To implement such a 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.

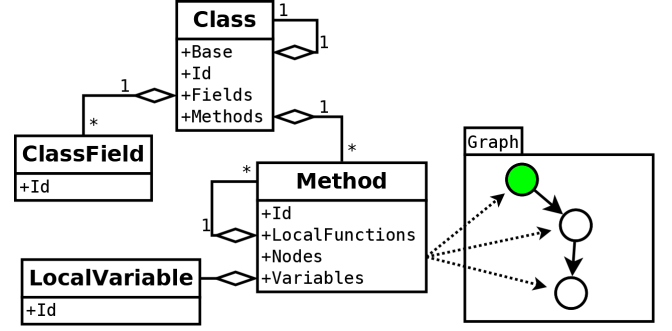


Figure 3: Structure of the metadata for extracted graph

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 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 trying to find all suitable entities.

2.2 Push-down automata construction

Further, we need to define restrictions on paths in the graph in terms of pushdown automata. Formally, nondeterministic push-down automaton [4] is a tuple $(Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$, where Q, Σ and Γ 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 takes the current state, an input symbol, the top of the stack and computes a new state and the sequence of stack symbols which must replace the top one. We also add the following restriction on the 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. It does not affect the expressive power of the resulting abstraction because push of a sequence of symbols can be emulated using the chain of states connected by ϵ -transitions each of which pushes one symbol.

Next, we propose to take the set of all edges in the control flow graph as Σ . 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 δ such that $\delta : Q \times \Sigma \cup \{\epsilon\} \times \Gamma \rightarrow \mathcal{P}(Q \times \Gamma^* \times N \cup \{v\})$, where N is the set of nodes of the graph and v is the dummy value that means that there is no need to jump and PDA just goes to the next node.

As an example, we construct the PDA performing the 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 δ is defined by the case analysis (4).

$$\begin{aligned}
1) \delta(q_0, i@invocation, \gamma) &:= \\
&\{(q_0, i\gamma, s_0), \dots, (q_0, i\gamma, s_n), (q_0, \gamma, v) : \\
&\quad s_0, \dots, s_n \in R(i)\} \\
2) \delta(q_0, a@assignment(v_s, v_t), \gamma) &:= \\
&\begin{cases} \{(v_t, \gamma, v), (q_0, \gamma, v)\}, & \text{if } source(v_s) \\ \{(q_0, \gamma, v)\}, & \text{otherwise} \end{cases} \\
3) \delta(v, a@assignment(v, v_t), \gamma) &:= (v_t, \gamma, v) \\
4) \delta(v, i@invocation, \gamma) &:= \\
&\bigcup_{j=0}^n \{(v_{j0}, i\gamma, s_j), \dots, (v_{jm}, i\gamma, s_j), (v, \gamma, v) : \\
&\quad v_{jk} \in A(i, j, v), s_j \in R(i)\} \\
5) \delta(v, r@return, i@invocation) &:= \\
&\begin{cases} \{(RV(i), \epsilon, T(i))\}, & \text{if } returned(v) \\ \emptyset, & \text{otherwise} \end{cases} \\
6) \delta(q, _, \gamma) &:= \{(q, \gamma, v)\}
\end{aligned} \tag{4}$$

Where notation $v@pattern$ means that v must be an object which is constructed by *pattern*, *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 (5).

$$A(i, j, v) := \begin{cases} \{q_f\}, & \text{if } j\text{-th target of invocation } i \\ & \text{is sink and } v \text{ is its argument} \\ \{v_k : v \mapsto v_k\}, & \text{if } j\text{-th target of } i \\ & \text{is not filter} \\ \emptyset, & \text{otherwise} \end{cases} \tag{5}$$

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.

2.3 Analysis execution

In this section we describe how the constructed automaton can be used to find an issue in the sample source code. The goal is to find the sequence of operations which starts in the entry point and ends at the invocation which uses the data from an unfiltered source. Since the automaton accepts such sequences, we can simulate the switching of its configurations according to the input statements taken from the source graph. We use the following notation for the PDA steps description. $(q, \gamma_1 \dots \gamma_n, n)$ is the current configuration of the simulation where q is the current state, γ_j is a symbol on the stack and n is the current position in the input graph. $c_1 \xrightarrow[k]{r} c_2$ is the k -th transition which switches configuration c_1 to c_2 using rule r from equation (4). Local variables are written as $\langle \text{Name of the containing method} \rangle. \langle \text{Variable identifier} \rangle$.

Since the automaton is non-deterministic, it can produce a graph of configurations, so, in this example, we explore only the branch where the final state is reached.

According to our example, we start with the configuration $(q_0, Z_0, 1)$. Configurations which appear after the first step are produced by accepting the first rule to the current configuration and the input symbol located at the edge between nodes 1 and 2. First of them is the one corresponding to the performed invocation of the *Read* method, the invocation statement is pushed onto the stack and the position is changed to 7. Second of them is the branch where invocation is just skipped. Since the second configuration does not lead to an error, we continue with the first configuration $(q_0, i_1 Z_0, 7)$ where i_1 is the invocation statement. Next step changes the state to the local $\langle \text{result} \rangle$ variable according to the rule 2 because the right part of the assignment is a source. Configuration switches to $(\text{Read}. \langle \text{result} \rangle, i_1 Z_0, 8)$. Further step processes the return statement using rule 5 and performs two important actions. Firstly, it pops the invocation stored on the top of the stack and jumps to the return point. Secondly, there is performed a change of currently tracked variable from the local one of the method *Read* to the local variable d of the method *Process* because it stores a result of the invocation. So, next configuration is $(\text{Process}.d, Z_0, 2)$. Processing of the edge between nodes 2 and 3 uses rule 4 and produces only one branch which just skips the invocation because the target is a filter and there is no need to enter this method. However, the invocation of *Consume* is processed fairly and there is produced the configuration $(\text{Consume}.d, i_2 Z_0, 11)$ which, further, iterates until the *Sink* invocation and reaches the final state q_f .

The full chain of configurations is shown in equation 6.

$$\begin{aligned}
(q_0, Z_0, 1) &\xrightarrow{1} (q_0, i_1 Z_0, 7) \xrightarrow{2} \\
&(\text{Read}. \langle \text{result} \rangle, i_1 Z_0, 8) \xrightarrow{5} \\
&(\text{Process}.d, Z_0, 2) \xrightarrow{4} (\text{Process}.d, Z_0, 3) \xrightarrow{5} \\
&(\text{Process}.d, i_2 Z_0, 10) \xrightarrow{6} (\text{Consume}.d, i_2 Z_0, 11) \xrightarrow{7} \\
&(q_f, i_3 i_2 Z_0, \langle \text{Sink entry point} \rangle)
\end{aligned} \tag{6}$$

So, the path in the graph which is present in equation (7) is accepted by the automaton and contains an issue.

$$\begin{aligned}
&\text{invoke:} && \text{assign:} \\
&\text{this.Read} \rightarrow \text{local d} && (\text{local a}).\text{Source} \rightarrow \langle \text{result} \rangle \\
1 \text{ -----} && \xrightarrow{\quad} && 7 \text{ -----} \\
&& \text{invoke:} && \text{invoke:} \\
&& \text{B.Filter} \rightarrow \text{local f} && \text{this.Consume} \\
8 \text{ -----} && \xrightarrow{\quad} && 2 \text{ -----} && \xrightarrow{\quad} && 3 \text{ -----} \\
&& (d := \text{local d}) && (d := \text{local d}) \\
&& \text{invoke:} && \text{invoke:} \\
&& \text{C.ctor} \rightarrow \text{local c} && (\text{local c}).\text{Sink} \\
10 \text{ -----} && \xrightarrow{\quad} && 11 \text{ -----} && \xrightarrow{\quad} && \langle \text{Sink entry point} \rangle
\end{aligned} \tag{7}$$

Where dashed arrows indicate jumps and solid ones correspond to straightforward transitions.

3 IMPLEMENTATION DETAILS

Using the described idea, we have developed the framework which makes it possible to implement interprocedural static code analysis based on CFL-reachability approach. Our 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 in fig. 4.

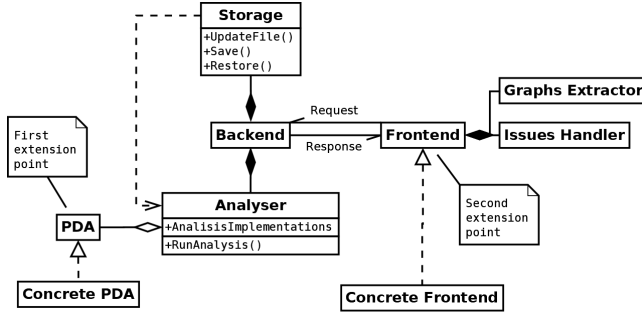


Figure 4: Solution structure

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 loading at the beginning of the next session. The second one is responsible for the execution of analyses. It contains an implementation of the resolver improving the one which is provided by IDE by adding dynamic invocations resolving such as lambdas propagation, 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 an implementation of the appropriate generic abstract class. Further, it is possible to run this new analysis using the existing internal algorithm of PDA simulation and get any finite subset of paths in the graph which are accepted by the PDA.

The second 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 one is a results interpreter which receives the set of paths in the graph each of which leads to an error, maps it to the source code and translates it to a human-readable format. For example, it can highlight pieces of code contained in the received paths.

Since a frontend is completely separate from the backend and the only requirement for it is to follow the communication protocol, the frontend can be considered as the second extension point. I.e. it is possible to replace the currently implemented frontend with any other implementation having the same functions as the original one including graphs extraction and results processing. The current implementation is also open to modifications which add support

for new types of analysis or any other features which requires interaction with IDE.

The protocol itself is based on request-response pattern 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. When there is a need to get the results of the analysis, for example, when the IDE performs the code highlighting, the frontend asks the backend for found issues, maps results to the source code and highlights corresponding lines.

4 EVALUATION

In order to test the resulting solution we have implemented the frontend as a plugin using ReSharper SDK², so it can be installed into ReSharper³, Rider⁴ and InspectCode⁵. The source code is parsed by internal ReSharper tools and the result is used to produce graphs and meta-information. The issues found by the backend are shown using code highlighting.

The sample analysis which has been implemented is the considered taint tracking analysis. It is defined by the PDA constructed in section 2 translated into the code with some slight modifications which make it possible to process interactions with object fields. To provide more information about an issue found by this analysis, the highlighting is accompanied by bulbs containing the full path of a tainted variable from the source to the sink represented as the sequence of operations.

4.1 Taint analysis sample cases

The resulting solution has been tested on a few common cases which can be found in the corresponding folder⁶. We take a closer look at three of them which help to illustrate the main properties of the analyzer. Each case is shown at a screenshot taken exactly from the run Rider IDE with some small relocations of bulbs to make them not to overlap the code.

Firstly, the solution ensures flow sensitivity. I.e. it processes the flow of variables passed into methods and returned from them distinguishing between different call sites and returning a result to the appropriate return point. It can be seen at fig 5. This example illustrates the most common cases of interprocedural data passing. *Brackets* method gets the data, possibly performs some computations on them and returns the result. Invocations at lines 37 and 38 show that the solution can distinguish two data flow paths despite both of them passes through the same method. So, *e* becomes tainted because *c* is tainted and *f* does not because *d* is clear. Moreover, the analyzer can track paths where passes and returns do not form the correct bracket sequence that is shown by method *PostSource* which does not take any parameter and just returns tainted data.

Secondly, the solution has limited context sensitivity. I.e. it allows tracking the propagation of objects that are tainted by assigning of

²ReSharper developer's guide: <https://www.jetbrains.com/help/resharper/sdk/README.html>. Access Date: 15.08.2019

³ReSharper: <https://www.jetbrains.com/resharper/>. Access Date: 15.08.2019

⁴Rider IDE: <https://www.jetbrains.com/rider/>. Access Date: 15.08.2019

⁵InspectCode tool: <https://www.jetbrains.com/help/resharper/InspectCode.html>. Access Date: 15.08.2019

⁶Test cases: github.com/gsvgit/CoFRA/tree/master/test/data/TaintAnalysis

```

17 class Program
18 {
19     [Tainted] private int A;
20     [Filter] private int Filter(int a) { return a; }
21     [Sink] private void Sink(int a) {}
22
23     private int PostSource() {
24         var b = A;
25         return b;
26     }
27
28     private int Brackets(int a) {
29         var b = a;
30         return b;
31     }
32
33     private static void Main(string[] args) {
34         var a = new Program();
35         var c = a.PostSource();
36         var d = a.Filter(c);
37         var e = a.Brackets(c);
38         var f = a.Brackets(d);
39         a.Sink(e);
40         a.Sink(f);
41     }
42 }

```

Tainted sink
source - Program.cs:24
assign - Program.cs:25
return -<-
assign - Program.cs:35
pass -> Program.cs:37 (System.Int32)Brackets(System.Int32)
assign - Program.cs:29
assign - Program.cs:30
return -<-
assign - Program.cs:37
sink -> Program.cs:39 (System.Void)Sink(System.Int32)

Figure 5: Flow sensitivity

some fields inside them both by their own methods and by outer code interacting with their fields directly. However, there is no possibility to provide true context-sensitivity since it cannot be expressed alongside with flow-sensitivity [10]. The first case is shown at fig 6. There is the field *B* at line 18. This field can be used widely in the logic of the *Container* class and by this, the tainting of this field is considered as the tainting of the whole object. However, while processing of the method *Store* during the analysis it is hard to decide what the object need to be tainted because in the inner context of *Store* it is just *this* object. I.e. we must consider the calling context to make such a decision. So, the solution provides this opportunity which is shown by lines 33-36 where the first invocation of *Store* leads to the tainting of object *d* and the second invocation does not taint object *e*.

Finally, the solution works with any type of recursion and does not fall into infinite cycles. It can be seen at fig. 7. This snippet contains two mutually recursive methods which pass the data to each other. The solution checks all possible paths of passing even those which includes cyclic invocations and returns the passed variable to the point corresponding to the initial invocation.

4.2 Performance

It is also necessary to measure the performance of the resulting solution. To do it more precisely, we run another type of analysis which tracks propagation of all variables and by this explores any possible path in the graph. Thus, the time and space required for its execution may be consistent estimation of the efficiency of the solution.

```

17 class Container {
18     private int B;
19     public void Store(int a) { B = a; }
20 }
21
22 class Program {
23     [Tainted] private int A;
24     [Filter] private int Filter(int a) { return a; }
25     [Sink] private void Sink(Container c) {}
26
27     private static void Main(string[] args) {
28         var a = new Program();
29         var b = a.A;
30         var c = a.Filter(b);
31         var d = new Container();
32         var e = new Container();
33         d.Store(b);
34         e.Store(c);
35         a.Sink(d);
36         a.Sink(e);
37     }
38 }

```

Tainted sink
source - Program.cs:29
pass -> Program.cs:33 (System.Void)Store(System.Int32)
assign - Program.cs:19
return -<-
sink -> Program.cs:35 (System.Void)Sink(TaintedTrackingTests.Container)

Figure 6: Tainting of an object by its own method

```

17 class Program
18 {
19     [Tainted] private int A;
20     [Filter] private int Filter(int a) { return a; }
21     [Sink] private void Sink(int a) {}
22
23     private int Recursive1(int c, int d) {
24         var r = Recursive2(c - 1, d);
25         return r;
26     }
27
28     private int Recursive2(int c, int d) {
29         if (c == 0) return d;
30         var r = Recursive1(c, d);
31         return r;
32     }
33
34     private static void Main(string[] args) {
35         var a = new Program();
36         var b = a.A;
37         var c = a.Filter(b);
38         var d = a.Recursive1(c, 10, d);
39         var e = a.Recursive1(c, 10, d);
40         a.Sink(d);
41         a.Sink(e);
42     }
43 }

```

Tainted sink
source - Program.cs:38
pass -> Program.cs:38 (System.Int32)Recursive1(System.Int32, System.Int32)
pass -> Program.cs:24 (System.Int32)Recursive2(System.Int32, System.Int32)
assign - Program.cs:23
return -<-
assign - Program.cs:24
assign - Program.cs:25
return -<-
assign - Program.cs:38
sink -> Program.cs:40 (System.Void)Sink(System.Int32)

Figure 7: Recursive methods processing

The code base which has been chosen as a source of data is the full solutions of a few big projects: Mono⁷, EventStore⁸ and OpenRA⁹. The analyzer has been tested on a computer running Windows 10 with quad-core Intel Core i7 3.4 GHz CPU and 16 GB of RAM. The results is shown in the table 1. Execution time does not include the time required for graph construction.

So, our solution is able to process real-world projects in acceptable time.

⁷Source code of mono project: <https://github.com/mono/mono>. Access Date: 15.08.2019

⁸Source code of EventStore project: <https://github.com/EventStore/EventStore>. Access Date: 15.08.2019

⁹Source code of OpenRA project: <https://github.com/OpenRA/OpenRA>. Access Date: 15.08.2019

Project	Classes	Methods	Execution time (s)	Allocated memory (GB)
Mono	21013	192745	21 ± 0.5	~ 4.2
EventStore	3828	22796	2.7 ± 0.1	~ 0.49
OpenRA	2767	15451	2.3 ± 0.1	~ 0.4

Table 1: Performance

5 CONCLUSION

We propose and implement in C# programming language the generic framework for interprocedural static code analysis implementation. This framework allows one to implement arbitrary interprocedural analysis in terms of CFL-reachability. By using the proposed framework, we implement a plugin upon ReSharper infrastructure which provides simple taint analysis and demonstrate that our solution can handle important real-world cases. Also we show that the proposed framework can be used for real-world solutions analysis.

One of the directions for future work is a creation of analysis and its evaluation on real-world projects. By this way, we want to get information which helps to improve the usability of our framework: tune performance, improve API, etc. Also we should improve documentation and create more examples of usage.

Another direction is a practical evaluation of automatic fix location prediction by using minimum cuts method [1].

Also we want to compare the proposed approach with other generic CFL-reachability based approaches for interprocedural code analysis creation. For example, with generation-based approach [3], which idea is similar to parser generators.

REFERENCES

- [1] Andrei Marian Dan, Manu Sridharan, Satish Chandra, Jean-Baptiste Jeannin, and Martin Vechev. 2017. Finding Fix Locations for CFL-Reachability Analyses via Minimum Cuts. In *Computer Aided Verification*, Rupak Majumdar and Viktor Kunčák (Eds.). Springer International Publishing, Cham, 521–541.
- [2] Jonas De Bleser, Quentin Stiévenart, Jens Nicolay, and Coen De Roover. 2017. Static Taint Analysis of Event-driven Scheme Programs. In *Proceedings of the 10th European Lisp Symposium*. ACM, 80–87.
- [3] Nicholas Hollingum and Bernhard Scholz. 2017. Cauliflower: a Solver Generator for Context-Free Language Reachability. In *LPAR-21. 21st International Conference on Logic for Programming, Artificial Intelligence and Reasoning (EPIC Series in Computing)*, Thomas Eiter and David Sands (Eds.), Vol. 46. EasyChair, 171–180. <https://doi.org/10.29007/tbm7>
- [4] John E. Hopcroft and Jeffrey D. Ullman. 1990. *Introduction To Automata Theory, Languages, And Computation* (1st ed.). Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA.
- [5] Wei Huang, Yao Dong, Ana Milanova, and Julian Dolby. 2015. Scalable and Precise Taint Analysis for Android. In *Proceedings of the 2015 International Symposium on Software Testing and Analysis (ISSTA 2015)*. ACM, New York, NY, USA, 106–117. <https://doi.org/10.1145/2771783.2771803>
- [6] Polyvios Pratikakis, Jeffrey S. Foster, and Michael Hicks. 2006. Existential Label Flow Inference Via CFL Reachability. In *Static Analysis*, Kwangkeun Yi (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 88–106.
- [7] Jakob Rehof and Manuel Fähndrich. 2001. Type-base Flow Analysis: From Polymorphic Subtyping to CFL-reachability. *SIGPLAN Not.* 36, 3 (Jan. 2001), 54–66. <https://doi.org/10.1145/373243.360208>
- [8] Thomas Reps. 1997. Program Analysis via Graph Reachability. In *Proceedings of the 1997 International Symposium on Logic Programming (ILPS '97)*. MIT Press, Cambridge, MA, USA, 5–19. <http://dl.acm.org/citation.cfm?id=271338.271343>
- [9] Thomas Reps. 2000. Undecidability of Context-sensitive Data-dependence Analysis. *ACM Trans. Program. Lang. Syst.* 22, 1 (Jan. 2000), 162–186. <https://doi.org/10.1145/345099.345137>
- [10] Johannes Späth, Karim Ali, and Eric Bodden. 2019. Context-, Flow-, and Field-sensitive Data-flow Analysis Using Synchronized Pushdown Systems. *Proc. ACM Program. Lang.* 3, POPL, Article 48 (Jan. 2019), 29 pages. <https://doi.org/10.1145/3290361>
- [11] Hao Tang, Di Wang, Yingfei Xiong, Lingming Zhang, Xiaoyin Wang, and Lu Zhang. 2017. Conditional Dyck-CFL Reachability Analysis for Complete and Efficient Library Summarization. In *Programming Languages and Systems*, Hongseok Yang (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 880–908.
- [12] Kai Wang, Aftab Hussain, Zhiqiang Zuo, Guoqing Xu, and Ardan Amiri Sani. 2017. Grasp: A Single-machine Disk-based Graph System for Interprocedural Static Analyses of Large-scale Systems Code. *SIGOPS Oper. Syst. Rev.* 51, 2 (April 2017), 389–404. <https://doi.org/10.1145/3093315.3037744>
- [13] Dacong Yan, Guoqing Xu, and Atanas Rountev. 2011. Demand-driven Context-sensitive Alias Analysis for Java. In *Proceedings of the 2011 International Symposium on Software Testing and Analysis (ISSTA '11)*. ACM, New York, NY, USA, 155–165. <https://doi.org/10.1145/2001420.2001440>
- [14] Qirun Zhang and Zhendong Su. 2017. Context-sensitive Data-dependence Analysis via Linear Conjunctive Language Reachability. *SIGPLAN Not.* 52, 1 (Jan. 2017), 344–358. <https://doi.org/10.1145/3093333.3009848>
- [15] Xin Zheng and Radu Rugina. 2008. Demand-driven Alias Analysis for C. In *Proceedings of the 35th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '08)*. ACM, New York, NY, USA, 197–208. <https://doi.org/10.1145/1328438.1328464>