

# Java Code Analysis with Database and Domain Specific Language

[www.huawei.com](http://www.huawei.com)

**Presenter:** Pan Linjie or 潘临杰 or денис

**Organization:** Huawei Cloud

**Date:** 2023-10-09



# Personal Profile



- Name: Pan Linjie or 潘临杰 or денис
- Organization: Huawei Cloud
- Work place: Beijing, China
- Education:
  - 2009—2013 Beihang University (BUAA) Bachelor
  - 2013—2016 Beihang University (BUAA) Master
  - 2016—2021 Institute of Software, Chinese Academy of Sciences (ISCAS) Doctor
- Research interests: Static Analysis, Software Testing, Constraint Solving
- Personal page: <https://github.com/pangeneral>

# Content

**01**

Background

**02**

Examples of the community

**03**

What we are doing

# **How to Perform Static Analysis?**

**Start from an Example**

# An Example of check rule

- 1. Define a check rule: *Case* should end with `break` or `$FALL-THROUGH$` comment
- 2. Implement the rule: write a checker
- 3. Scan the code

## Right case

```
switch (label) {
    case 0:
    case 1:
        System.out.println("1");
        // $FALL-THROUGH$
    case 2:
        System.out.println("2");
        // $FALL-THROUGH$
    case 3:
        System.out.println("3");
        break;
    default:
        System.out.println("Default case!");
}
```

## Wrong case

```
switch (label) {
    case 0:
    case 1:
        System.out.println("1");
    case 2:
        System.out.println("2");
    case 3:
        System.out.println("3");
        break;
    default:
        System.out.println("Default case!");
}
```

# An Example about checker

- 1. Define a check rule: *Case* should end with *break* or *\$FALL-THROUGH\$* comment
- 2. Implement the rule: write a checker
- 3. Scan the code

```
1 public class FallThroughFixer extends JavaGenericDefectFixer {
2     public List<DefectInstance> detectDefectsForFileModel(JavaClass javaClass) {
3         CompilationUnit compilationUnit = javaClass.getCompilationUnit();
4         compilationUnit.accept(new ASTVisitor() {
5             @Override
6             public boolean visit(SwitchStatement node) {
7                 checkFallThrough(node.statements(), javaClass);
8                 return super.visit(node);
9             }
10        });
11        return getCurrentFileDefectInstances();
12    }
13
14    private void checkFallThrough(List<ASTNode> switchStmts, JavaClass javaClass) {
15        Map<SwitchCase, List<ASTNode>> caseGroup = caseGroup(switchStmts);
16        CompilationUnit compilationUnit = javaClass.getCompilationUnit();
17        SwitchCase lastCase = caseGroup.keySet().stream()
18            .max(Comparator.comparingInt(node -> compilationUnit.getLineNumber(node.getStartPosition())))
19            .orElse(null);
20        caseGroup.forEach((switchCase, nodes) -> {
21            if (CollectionUtils.isEmpty(nodes)) {
22                return;
23            }
24            boolean hasBreak = nodes.stream().anyMatch(this::hasBreak);
25            if (!hasBreak && !switchCase.equals(lastCase)) {
26                ASTNode lastNode = nodes.get(nodes.size() - 1);
27                int endLineNum = compilationUnit.getLineNumber(lastNode.getStartPosition() + lastNode.getLength());
28                String potentialCmt = getCurrentFileLines().get(endLineNum);
29                if (potentialCmt.trim().startsWith("//") || potentialCmt.trim().startsWith("/*")) {
30                    return;
31                }
32                int bugLineNum = compilationUnit.getLineNumber(switchCase.getStartPosition());
33                int bugColumnNum = compilationUnit.getColumnNumber(switchCase.getStartPosition()) + 1;
34                DefectInstance defectInstance = createADummyDefectInstance(bugLineNum, bugLineNum, bugColumnNum);
35                defectInstance.getContext().add("lastNode", lastNode);
36                defectInstance.getContext().add("javaClass", javaClass);
37                saveDefectInstance(defectInstance);
38            }
39        });
40    }
```

```
42     private boolean hasBreak(ASTNode node) {
43         AtomicBoolean hasBreak = new AtomicBoolean(false);
44         node.accept(new ASTVisitor() {
45             @Override
46             public boolean visit(BreakStatement node) {
47                 hasBreak.set(true);
48                 return super.visit(node);
49             }
50         });
51         return hasBreak.get();
52     }
53
54     private Map<SwitchCase, List<ASTNode>> caseGroup(List<ASTNode> switchStmts) {
55         Map<SwitchCase, List<ASTNode>> caseGroup = new HashMap<>();
56         for (int i = 0; i < switchStmts.size(); i++) {
57             ASTNode stt = switchStmts.get(i);
58             if (!(stt instanceof SwitchCase)) {
59                 continue;
60             }
61             SwitchCase switchCase = (SwitchCase) stt;
62             int index = i + 1;
63             List<ASTNode> list = new ArrayList<>();
64             ASTNode nextNode = null;
65             while (!(nextNode instanceof SwitchCase) && index < switchStmts.size()) {
66                 nextNode = switchStmts.get(index);
67                 if (!(nextNode instanceof SwitchCase)) {
68                     list.add(nextNode);
69                 }
70                 index++;
71             }
72             caseGroup.put(switchCase, list);
73         }
74         return caseGroup;
75     }
```

1 CommentIndentationTest.java

G.FMT.16 case语句块结束时如果不加break, 需要有注释说明(fall-through)

提示 待处理 责任人: [redacted] 修改问题状态 申请屏蔽 问题帮助 查看上下文

```
31 switch (obj) {
32     case "":
33         System.out.println();
34         default;
```

2 FallThroughTest.java

G.FMT.16 case语句块结束时如果不加break, 需要有注释说明(fall-through)

提示 待处理 责任人: [redacted] 修改问题状态 申请屏蔽 问题帮助 查看上下文

```
22 // CHECK-FIXLINES @@LINE+1]$:line: break [FallThrough]
23 System.out.println();
24 case "":
25 // CHECK-MESSAGES @@LINE-1] 13 warning if no break is added at the end of the case statement block, a comment [fall-through] is required [FallThrough]
26 // CHECK-FIXLINES @@LINE+5]$:line: [FallThrough]
```

3 FallThroughTest.java

G.FMT.16 case语句块结束时如果不加break, 需要有注释说明(fall-through)

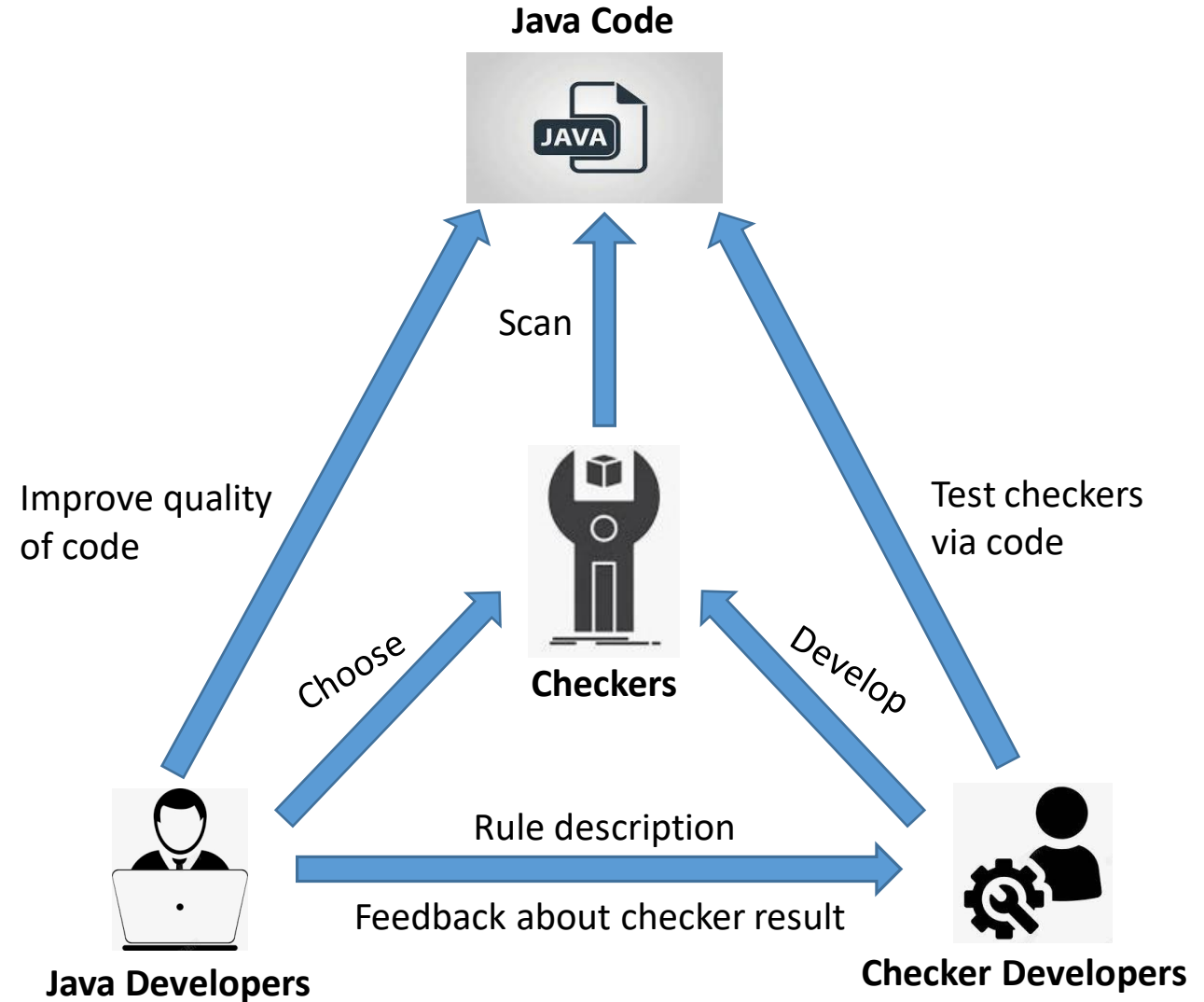
提示 待处理 责任人: [redacted] 修改问题状态 申请屏蔽 问题帮助 查看上下文

```
17 // fall through
18 case "":
19 case "":
20 // CHECK-MESSAGES @@LINE-1] 13 warning if no break is added at the end of the case statement block, a comment [fall-through] is required [FallThrough]
21 // CHECK-FIXLINES @@LINE+2]$:line: System.out.println(); [FallThrough]
```

1. Relatively complex checker logic compared with simple rule description
2. Operation on nodes of abstract syntax tree: *JavaClass*, *ASTNode*, *SwitchStatement*

# Static Analysis in Industry

- **Targets** of analysis checkers: Java code
  - Huge scale
  - Complex logic
  - Harmful bugs and vulnerabilities
- **Users** of analysis checkers: Java developers
  - Fast iterative development
  - Locate the potential bugs and vulnerabilities
  - Quick feedback about the result of static analysis
- **Owners** of analysis checkers: checker developers
  - Various requirements
  - Customized check rule
  - Develop and modify check rules timely



**Simplify the development of checkers**

# Domain-Specific Language (DSL)

- General-Purpose Language (GPL): JAVA, C++
- Common DSL: Maven, Gradle, HTML, SQL
- Our domain: static analysis
  - Program is very complex
    - Rice theorem: any non-trivial semantic property of a language which is recognized by a Turing machine is **undecidable**
    - The scene of analysis is various
  - Different scenes are similar
    - Fixed program elements: class, method, statement, expression, variable
    - Obtain program elements and perform analysis logic
- The goal of DSL-based analysis
  - Simplify development of analysis rule
  - Reduce duplicate analysis

```
1 public class FallThroughFixer extends JavaGenericDefectFixer {
2     public List<DefectInstance> detectDefectsForFileModel(JavaClass javaClass) {
3         CompilationUnit compilationUnit = javaClass.getCompilationUnit();
4         compilationUnit.accept(new ASTVisitor() {
5             @Override
6             public boolean visit(SwitchStatement node) {
7                 checkFallThrough(node.statements(), javaClass);
8                 return super.visit(node);
9             }
10        });
11        return getCurrentFileDefectInstances();
12    }
13
14    private void checkFallThrough(List<ASTNode> switchStmts, JavaClass javaClass) {
15        Map<SwitchCase, List<ASTNode>> caseGroup = caseGroup(switchStmts);
16        CompilationUnit compilationUnit = javaClass.getCompilationUnit();
17        SwitchCase lastCase = caseGroup.keySet().stream()
18            .max(Comparator.comparingInt(node -> compilationUnit.getLineNumber(node.getStartPosition())))
19            .orElse(null);
20        caseGroup.forEach((switchCase, nodes) -> {
21            if (CollectionUtils.isEmpty(nodes)) {
22                return;
23            }
24            boolean hasBreak = nodes.stream().anyMatch(this::hasBreak);
25            if (!hasBreak && !switchCase.equals(lastCase)) {
26                ASTNode lastNode = nodes.get(nodes.size() - 1);
27                int endLineNum = compilationUnit.getLineNumber(lastNode.getStartPosition() + lastNode.getLength());
28                String potentialCmt = getCurrentFileLines().get(endLineNum);
29                if (potentialCmt.trim().startsWith("//") || potentialCmt.trim().startsWith("/*")) {
30                    return;
31                }
32                int bugLineNum = compilationUnit.getLineNumber(switchCase.getStartPosition());
33                int bugColumnNum = compilationUnit.getColumnNumber(switchCase.getStartPosition() + 1);
34                DefectInstance defectInstance = createADummyDefectInstance(bugLineNum, bugLineNum, bugColumnNum);
35                defectInstance.getContext().add("lastNode", lastNode);
36                defectInstance.getContext().add("javaClass", javaClass);
37                saveDefectInstance(defectInstance);
38            }
39        });
40    }
41 }
```



Find caseBlock cb where  
cb.lastStatement is not or(break, fallthrough)

## How to represent code?



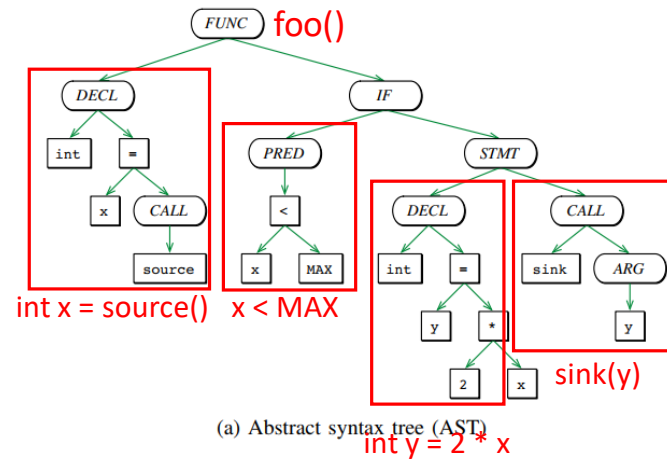
# Abstraction of Code

```

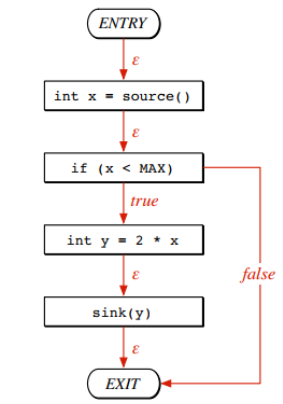
void foo()
{
  int x = source();
  if (x < MAX)
  {
    int y = 2 * x;
    sink(y);
  }
}

```

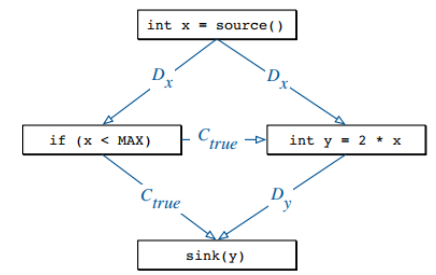
1  
2  
3  
4  
5  
6  
7  
8  
9



(a) Abstract syntax tree (AST)



(b) Control flow graph (CFG)



(c) Program dependence graph (PDG)

- Representation of code

- Abstract syntax tree: syntactic structure
- Control flow graph: execution path
- Program dependence graph: data and control dependence

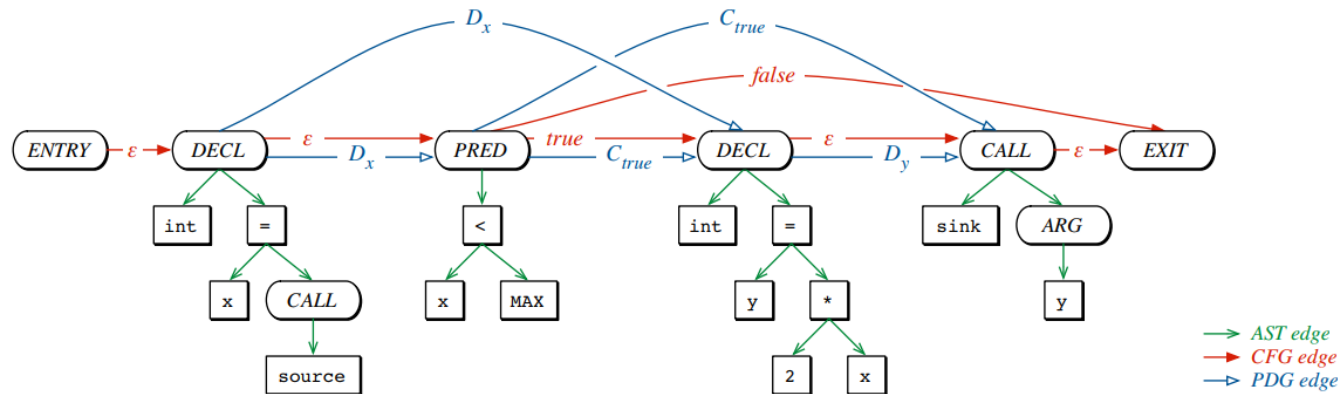
# About Code Property Graph

- Code property graph: the combination of AST, CFG and PDG [1]

**Definition 3.** A code property graph is a property graph  $G = (V, E, \lambda, \mu)$  constructed from the AST, CFG and PDG of source code with

- $V = V_A$ ,
- $E = E_A \cup E_C \cup E_P$ ,
- $\lambda = \lambda_A \cup \lambda_C \cup \lambda_P$  and
- $\mu = \mu_A \cup \mu_E$ ,

where we combine the labeling and property functions with a slight abuse of notation.



Gather information of code together

Extend CPG to contain more information

## Modeling and Discovering Vulnerabilities with Code Property Graphs

Fabian Yamaguchi\*, Nico Golde<sup>†</sup>, Daniel Arp\* and Konrad Rieck\*  
<sup>\*</sup>University of Göttingen, Germany  
<sup>†</sup>Qualcomm Research Germany

**Abstract**—The vast majority of security breaches encountered today are a direct result of insecure code. Consequently, the protection of computer systems critically depends on the rigorous identification of vulnerabilities in software, a tedious and error-prone process requiring significant expertise. Unfortunately, a single flaw suffices to undermine the security of a system and thus the sheer amount of code to audit plays into the attacker's cards. In this paper, we present a method to effectively mine large amounts of source code for vulnerabilities. To this end, we introduce a novel representation of source code called a code property graph that merges concepts of classic program analysis, namely abstract syntax trees, control flow graphs and program dependence graphs, into a joint data structure. This comprehensive representation enables us to elegantly model templates for common vulnerabilities with graph traversals that, for instance, can identify buffer overflows, integer overflows, format string vulnerabilities, or memory disclosures. We implement our approach using a popular graph database and demonstrate its efficacy by identifying 18 previously unknown vulnerabilities in the source code of the Linux kernel.

**Keywords**—Vulnerabilities; Static Analysis; Graph Databases

### I. INTRODUCTION

The security of computer systems fundamentally depends on the quality of its underlying software. Despite a long series of research in academia and industry, security vulnerabilities regularly manifest in program code, for example as failures to account for buffer boundaries or as insufficient validation of input data. Consequently, vulnerabilities in software remain one of the primary causes for security breaches today. For example, in 2013 a single buffer overflow in a universal plug-and-play library rendered over 23 million routers vulnerable to attacks from the Internet [26]. Similarly, thousands of users currently fall victim to web-based malware that exploits different flaws in the Java runtime environment [29].

The discovery of software vulnerabilities is a classic yet challenging problem of security. Due to the inability of a program to identify non-trivial properties of another program, the generic problem of finding software vulnerabilities is undecidable [33]. As a consequence, current means for spotting security flaws are either limited to specific types of vulnerabilities or build on tedious and manual auditing. In particular, securing large software projects, such as an operating system kernel, resembles a daunting task, as a single flaw may undermine the security of the entire code base. Although some classes of vulnerabilities recurring throughout the software landscape exist for a long time, such as buffer overflows and format string vulnerabilities, automatically detecting their

incarnations in specific software projects is often still not possible without significant expert knowledge [16].

As a result of this situation, security research has initially focused on statically finding specific types of vulnerabilities, such as flaws induced by insecure library functions [6], buffer overflows [45], integer overflows [40] or insufficient validation of input data [18]. Based on concepts from software testing, a broader detection of vulnerabilities has then been achieved using dynamic program analysis, ranging from simple fuzz testing [e.g., 38, 42] to advanced taint tracking and symbolic execution [e.g., 2, 35]. While these approaches can discover different types of flaws, they are hard to operate efficiently in practice and often fail to provide appropriate results due to either prohibitive runtime or the exponential growth of execution paths to consider [16, 21]. As a remedy, security research has recently started to explore approaches that assist an analyst during auditing instead of replacing her. The proposed methods accelerate the auditing process by augmenting static program analysis with expert knowledge and can thereby guide the search for vulnerabilities [e.g., 39, 43, 44].

In this paper, we continue this direction of research and present a novel approach for mining large amounts of source code for vulnerabilities. Our approach combines classic concepts of program analysis with recent developments in the field of graph mining. The key insight underlying our approach is that many vulnerabilities can only be adequately discovered by jointly taking into account the structure, control flow and dependencies of code. We address this requirement by introducing a novel representation of source code denoted as code property graph. This graph combines properties of abstract syntax trees, control flow graphs and program dependence graphs in a joint data structure. This comprehensive view on code enables us to elegantly model templates for common vulnerabilities using graph traversals. Similar to the query in a database, a graph traversal passes over the code property graph and inspects the code structure, the control flow, and the data dependencies associated with each node. This joint access to different code properties enables crafting concise templates for several types of flaws and thereby helps to audit large amounts of code for vulnerabilities.

We implement our approach using a popular graph database and demonstrate its practical merits by designing graph traversals for several well-known vulnerability types, such as buffer overflows, integer overflows, format string vulnerabilities, or memory disclosures. As a show case, we analyze the source code of the Linux kernel—a large and well-audited code base. We find that almost all vulnerabilities reported for

## How to save code?

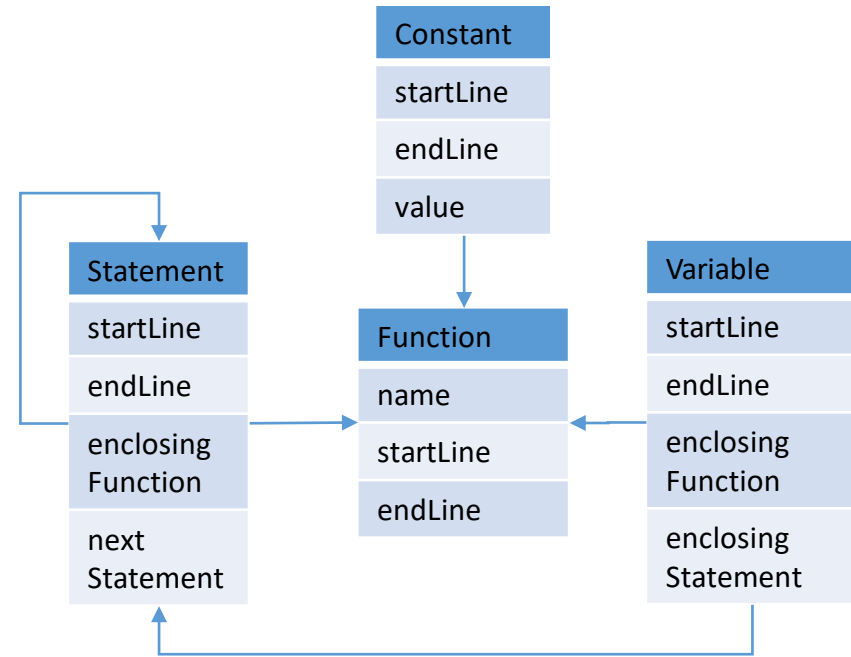
# Relational Database VS Graph Database

```

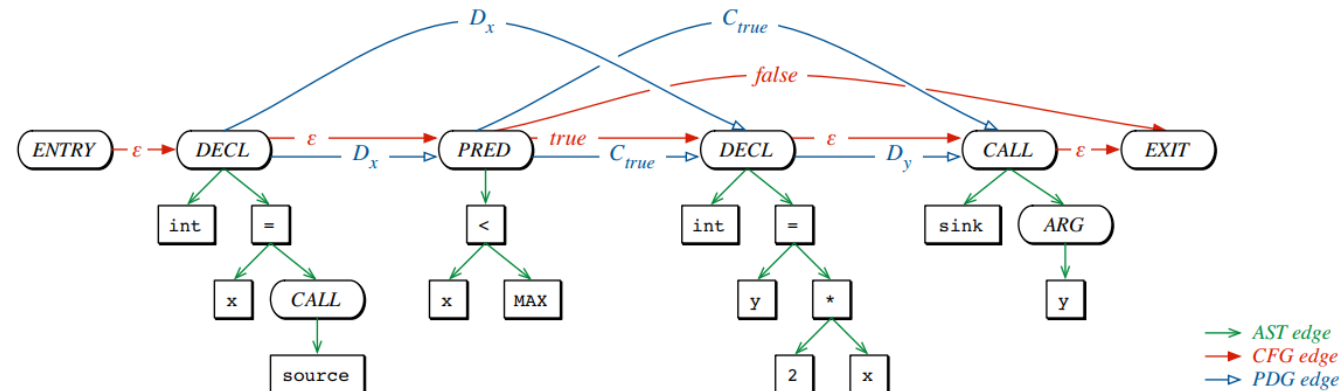
void foo()
{
    int x = source();
    if (x < MAX)
    {
        int y = 2 * x;
        sink(y);
    }
}

```

1  
2  
3  
4  
5  
6  
7  
8  
9



	Relational Database	Graph Database
Data organization	Table and relation	Node and edge
Strength	Filter structured data	Process relation



- Which one is better?


How many classes do we have?

Is there a path from a statement to another?

**Relation is more important than data**

# Common Graph Database

Commercialization



- **TinkerGraph:** light、 easy to configure、 support Gremlin、 no transaction、 open-source
- **OverflowDB:** memory friendly、 open-source
- **JanusGraph:** support huge graph、 concurrent transaction、 open-source of code
- **TigerGraph:** concurrent graph algorithm big data process、 license for commercial application
- **Amazon Neptune:** High performance、 multiple levels of security、 Must access via virtual private cloud
- **Neo4j:** master-slave topology cluster、 real-time schema update、 license for commercial application

# Content

**01**

Background

**02**

Examples of the community

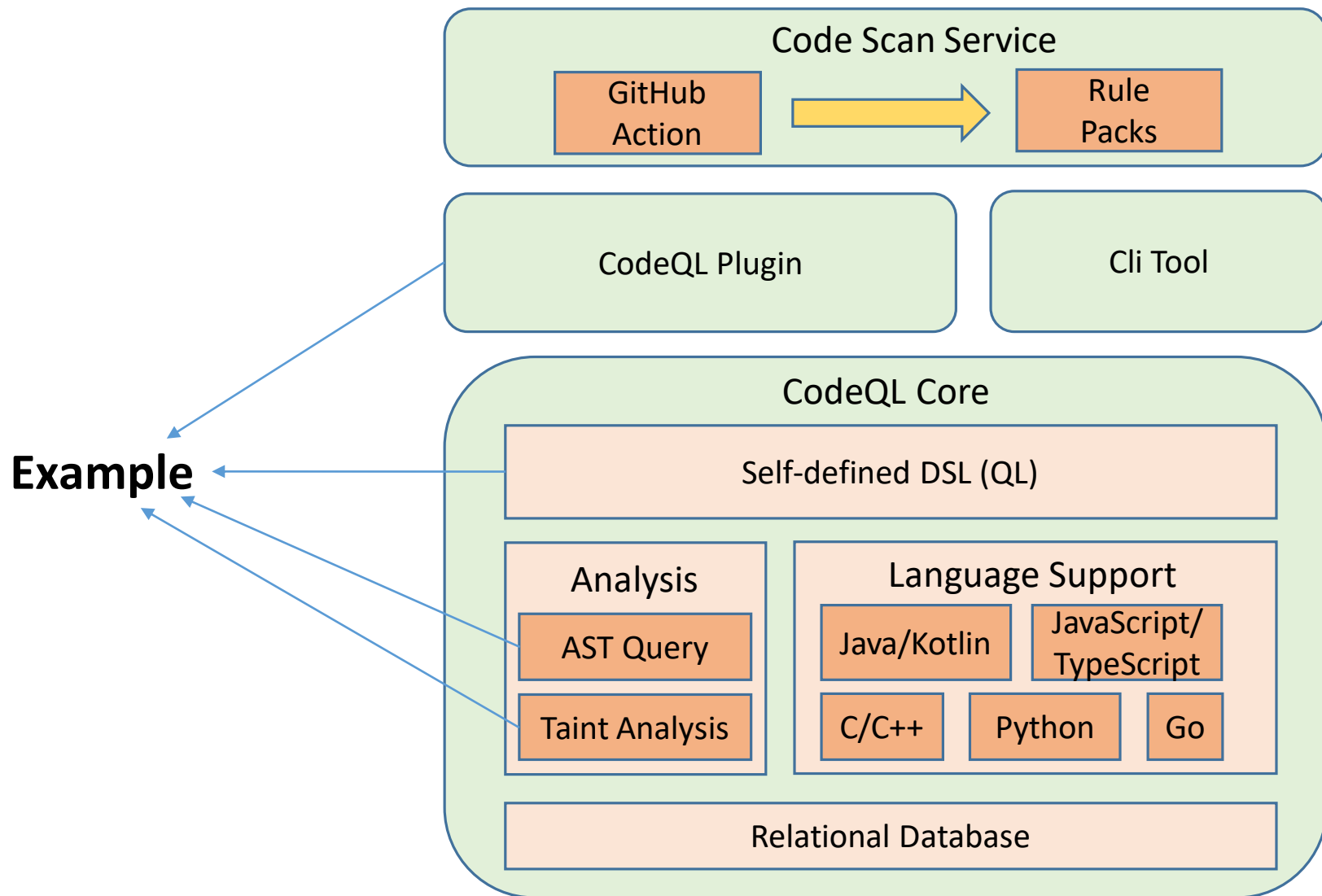
**03**

What we are doing

# Tools using Database and DSL

- CodeQL: relational database with self-defined DSL
- Joern: graph database with Gremlin
- Codyze: graph database with Gremlin

# Framework of CodeQL



# An Example of Query Analysis via CodeQL

- Rule: find all empty blocks
- Structure of checker
  - Metadata
  - Import CodeQL libraries
  - CodeQL class and predicates definition
  - CodeQL variable declaration
  - CodeQL query logic
  - Entrance of query: select clause
- Query kind
  - Alert (@kind problem): highlight issues in specific code locations
  - Path (@kind path-problem): describe the flow of information between a source and a sink

Not Java, Not Java, Not Java

```
1  /**
2   * @name Empty block
3   * @kind problem
4   * @problem.severity warning
5   * @id java/example/empty-block
6   */
7
8  import java
9  import semmle.code.java.dataflow.DataFlow
10
11 from BlockStmt b
12
13 where b.getNumStmt() = 0
14
15 select b, "This is an empty block."
```

Self-defined DSL called QL



# An Example of Taint Analysis via CodeQL

- Import dependencies
- Inherit TaintTracking::Configuration
- Define source and sink
  - Source: string contant "tainted"
  - Sink: the first argument of function call
- Define taint condition
- Code under analysis

```
27 public void dataflowTest() {
28     String x = "tainted";
29     String y = "still" + x;
30     System.out.println(y);
31 }
32 }
33
```

- Analysis result

```
alerts 1 result Show results in Problems view
Message
<message> DataFlowTest.java:30:28
Path
1 "tainted" : String DataFlowTest.java:28:20
2 y DataFlowTest.java:30:28
```

```
/**
 * @kind path-problem
 */

import java
import semmler.code.java.dataflow.DataFlow
import semmler.code.java.dataflow.TaintTracking
import DataFlow::PathGraph

class TaintConf extends TaintTracking::Configuration {
    Quick Evaluation: TaintConf
    TaintConf() {
        this = "Taint Configuration"
    }

    Quick Evaluation: isSource
    override predicate isSource(DataFlow::Node source) {
        source.asExpr().(StringLiteral).getValue() = "tainted"
    }

    Quick Evaluation: isSink
    override predicate isSink(DataFlow::Node sink) {
        exists(Call call |
            sink.asExpr() = call.getArgument(0)
        )
    }
}

from TaintConf config,
DataFlow::PathNode source, DataFlow::PathNode sink
where config.hasFlowPath(source, sink)
select sink.getNode(), source, sink, "<message>"
```

# Taint Analysis with Sanitizer and Additional Flow

- Source: string constant "source"
- Sink: the first argument of function call println
- Sanitizer: the first argument of function call sanitize
- Self-defined flow: from the first argument to the second of function call extraFlow

```
28 public void dataFlowTest() {
29     // source
30     String source = "source";
31     String directTainted = "direct" + source;
32     String taintedByExtraFlow = "ExtraFlow";
33     this.extraFlow(source, taintedByExtraFlow);
34     // sink point 1
35     System.out.println(taintedByExtraFlow);
36     this.sanitize(source);
37     String notTaintedByExtraFlow = "notTainted";
38     this.extraFlow(source, notTaintedByExtraFlow);
39     // sink point 2
40     System.out.println(directTainted);
41     System.out.println(notTaintedByExtraFlow); // not sink
42 }
43
44 1 usage
45 public void sanitize(String str) {
46 }
47
48 2 usages
49 public void extraFlow(String first, String second) {
50     System.out.printf("{} is tainted by {}", second, first);
51 }
```

```
class TaintConf extends TaintTracking::Configuration {
    Quick Evaluation: TaintConf
    TaintConf() {
        this = "Taint Configuration"
    }

    Quick Evaluation: isSource
    override predicate isSource(DataFlow::Node source) {
        source.asExpr().(StringLiteral).getValue() = "source"
    }

    Quick Evaluation: isSink
    override predicate isSink(DataFlow::Node sink) {
        exists(MethodAccess ma |
            sink.asExpr() = ma.getArgument(0) and
            ma.getMethod().getName() = "println"
        )
    }

    Quick Evaluation: isSanitizer
    override predicate isSanitizer(DataFlow::Node node) {
        exists(MethodAccess ma |
            ma.getArgument(0) = node.asExpr() and
            ma.getMethod().getName() = "sanitize"
        )
    }

    Quick Evaluation: isAdditionalTaintStep
    override predicate isAdditionalTaintStep(DataFlow::Node node1, DataFlow::Node node2) {
        exists(MethodAccess ma |
            ma.getArgument(0) = node1.asExpr() and
            ma.getArgument(1) = node2.asExpr() and
            ma.getMethod().getName() = "extraFlow"
        )
    }
}
```

alerts 2 results Show results in Problems view

Message	
<message>	DataFlowTest.java:35:28
Path	
1 "source": String	DataFlowTest.java:30:25
2 taintedByExtraFlow	DataFlowTest.java:35:28
<message>	DataFlowTest.java:40:28
Path	
1 "source": String	DataFlowTest.java:30:25
2 directTainted	DataFlowTest.java:40:28

# Taint Analysis with State

- Source
  - 1. string constant "sourceCleaned", initialized with state "cleaned"
  - 2. string constant "sourceNotCleaned", initialized with state "notCleaned"
- Sink: the first argument of function call println with state "cleaned", "notCleaned" or "taintFromCleaned"
- Sanitizer: the first argument of sanitize with state "cleaned"
- Self-defined flow: Taint the second argument, initialized with state "taintFromCleaned", if the first argument has the state "cleaned"

```
28 public void dataFlowTest() {
29     // sourceCleaned
30     String sourceCleaned = "sourceCleaned";
31     // sourceNotCleaned
32     String sourceNotCleaned = "sourceNotCleaned";
33     String tainted = "";
34     String notTainted = "";
35
36     // no additional flow
37     this.extraFlow(sourceNotCleaned, notTainted);
38     // additional flow
39     this.extraFlow(sourceCleaned, tainted);
40     this.sanitize(sourceNotCleaned);
41     this.sanitize(sourceCleaned);
42
43     System.out.println(sourceCleaned);
44     System.out.println(sourceNotCleaned); // sink point
45     System.out.println(notTainted);
46     System.out.println(tainted); // sink point
47 }
48
49 public void sanitize(String str) {}
50
51 public void extraFlow(String first, String second) {
52     System.out.printf("{} is tainted by {}", second, first);
53 }
54 }
```

Clean the taint

```
class TaintConf extends TaintTracking::Configuration {
    Quick Evaluation: TaintConf
    TaintConf() {
        this = "Taint Configuration"
    }
}

Quick Evaluation: isSource
override predicate isSource(DataFlow::Node source, DataFlow::FlowState state) {
    (source.asExpr().(StringLiteral).getValue() = "sourceCleaned" and state = "cleaned") or
    (source.asExpr().(StringLiteral).getValue() = "sourceNotCleaned" and state = "notCleaned")
}

Quick Evaluation: isSink
override predicate isSink(DataFlow::Node sink, DataFlow::FlowState state) {
    exists(MethodAccess ma |
        sink.asExpr() = ma.getArgument(0) and
        ma.getMethod().getName() = "println" and
        (state = "cleaned" or state = "notCleaned" or state = "taintFromCleaned")
    )
}

Quick Evaluation: isSanitizer
override predicate isSanitizer(DataFlow::Node node, DataFlow::FlowState state) {
    exists(MethodAccess ma |
        ma.getArgument(0) = node.asExpr() and
        ma.getMethod().getName() = "sanitize" and
        state = "cleaned"
    )
}

override predicate isAdditionalTaintStep(DataFlow::Node node1, DataFlow::FlowState state1,
    DataFlow::Node node2, DataFlow::FlowState state2) {
    exists(MethodAccess ma |
        ma.getArgument(0) = node1.asExpr() and
        ma.getArgument(1) = node2.asExpr() and
        ma.getMethod().getName() = "extraFlow" and
        state1 = "cleaned" and state2 = "taintFromCleaned")
    }
}
```

alerts 2 results Show results in Problems view

Message		
<message>		DataFlowTest.java:44:28
Path		
1	"sourceNotCleaned" : String	DataFlowTest.java:32:35
2	sourceNotCleaned	DataFlowTest.java:44:28
<message>		DataFlowTest.java:46:28
Path		
1	"sourceCleaned" : String	DataFlowTest.java:30:32
2	sourceCleaned : String	DataFlowTest.java:39:24
3	tainted : String	DataFlowTest.java:39:39
4	tainted	DataFlowTest.java:46:28

# CodeQL Plugin **only** available in VSCode

DSL Code (Checker)

View analysis result

Import database under analysis

Query history

View AST of code under analysis

The screenshot displays the Visual Studio Code interface with the CodeQL extension. The left sidebar contains the 'CODEQL' view with sections for 'DATABASES', 'VARIANT ANALYSIS REPOSITORIES', 'QUERY HISTORY', and 'AST VIEWER'. The main editor shows a CodeQL query file 'example.q1' and a Java source file 'AbstractExecutor.java'. The right sidebar shows the 'CodeQL Query Results' panel with a list of alerts, all of which are 'This is an empty block'. The bottom status bar shows the current file path and line/col information.

```
codeql-custom-queries-java > example.q1 > {} example
1 /**
2  * @name Empty block
3  * @kind problem
4  * @problem.severity warning
5  * @id java/example/empty-block
6  */
7
8 import java
9
10 from BlockStmt b
11 where b.getNumStmt() = 0
12 select b, "This is an empty block."
13 ^
```

```
AbstractExecutor.java
E:\lite-white-tiger\codeql\src.zip > E:\lite-white-tiger > secBrella-jvm-core > src > main > java > com > huawei > secBrella > jvm > core
51     LOGGER.info("executing time: [{}]", watch.elapsed(TimeUnit.SECONDS));
52     LOGGER.info("end [{}]", this.step());
53 }
54
55 /**
56  * 正在分析的入口
57  *
58  * @throws SecBrellaException 异常
59  */
60 protected abstract void run() throws SecBrellaException;
61
62 /**
63  * 分析完成后执行
64  *
65  * @throws SecBrellaException 异常
66  */
67 protected void postExecute() throws SecBrellaException {
68 }
69
70 }
```

CodeQL Query Results: Empty block on codeql - finished in 0 seconds (7 results) [4/7/2023, 10:36:55 AM]

Message	Location
This is an empty block.	AbstractExecutor.java:67:60
This is an empty block.	AbstractRenderer.java:54:67
This is an empty block.	AbstractRenderer.java:64:113
This is an empty block.	AbstractRenderer.java:73:65
This is an empty block.	AbstractConfigurableDetector.java:57:27
This is an empty block.	AbstractConfigurableJavaDetector.java:54:27
This is an empty block.	JavaParserUtil.java:428:83

PROBLEMS 8 OUTPUT TERMINAL DEBUG CONSOLE

```
Resolving tests using CodeQL CLI: resolve tests --strict-test-discovery -v --log-to-stderr --format json e:\Code\vscode-codeql-starter\codeql-custom-queries-java...
CLI command succeeded.

Resolving tests using CodeQL CLI: resolve tests --strict-test-discovery -v --log-to-stderr --format json e:\Code\vscode-codeql-starter\codeql-custom-queries-java...
CLI command succeeded.

Resolving tests using CodeQL CLI: resolve tests --strict-test-discovery -v --log-to-stderr --format json e:\Code\vscode-codeql-starter\codeql-custom-queries-java...
CLI command succeeded.

Resolving tests using CodeQL CLI: resolve tests --strict-test-discovery -v --log-to-stderr --format json e:\Code\vscode-codeql-starter\codeql-custom-queries-java...
CLI command succeeded.
```

# Content

**01**

Background

**02**

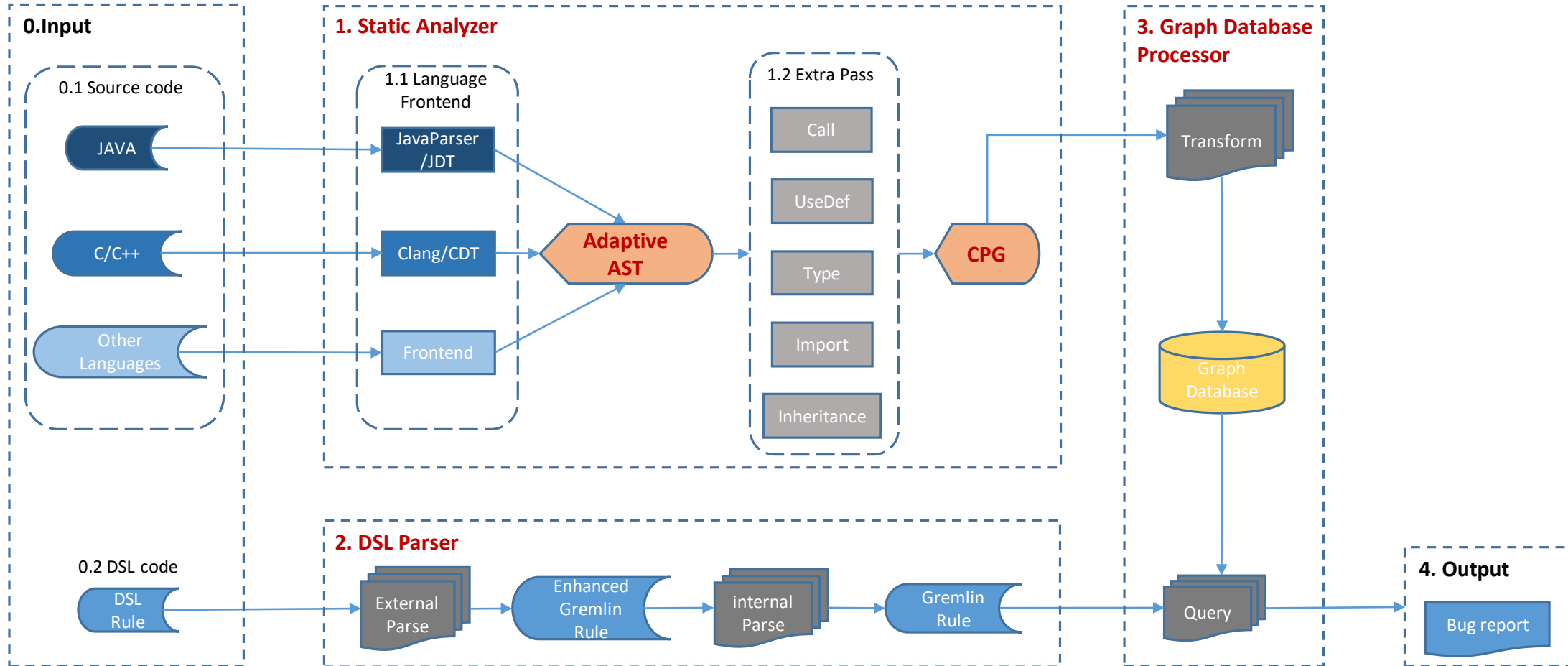
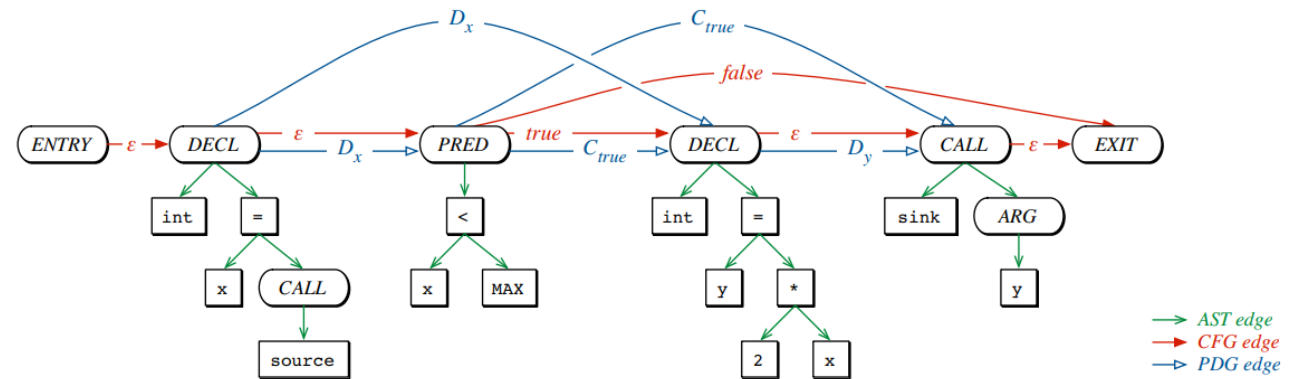
Examples of the community

**03**

What we are doing

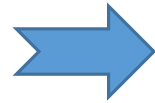
# Workflow of Our Analysis Engine

- Graph database with self-defined DSL
- Adaptive AST and CPG
- DSL translation



# Adaptive AST and CPG

- Java AST Schema
  - Class Declaration
  - Function Declaration
  - Variable Declaration
  - **Import Declaration**
  - ...
- C++ AST Schema
  - Class Declaration
  - Function Declaration
  - Variable Declaration
  - **Pointer Type**
  - ...



## Is it Unified Schema?

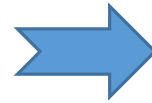
- Adaptive AST Schema
  - Class Declaration
  - Function Declaration
  - Variable Declaration
  - **Import Declaration**
  - **Pointer Type**
  - ...



Python  
Schema



Typescript  
Schema

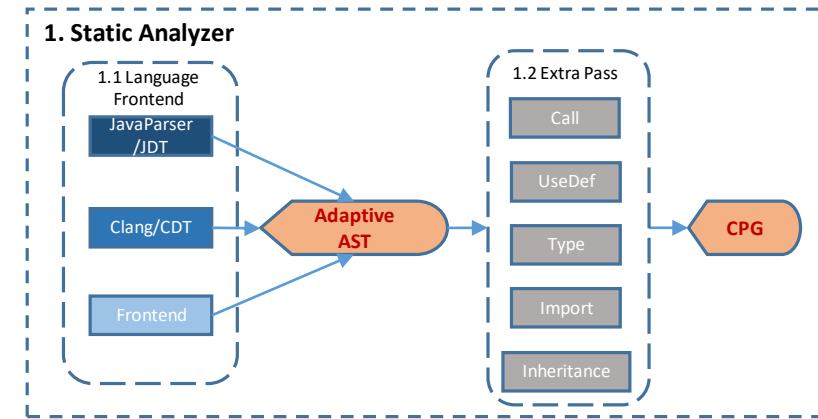


- Code Property Graph

- AST Nodes and edges
- **Control flow edges**
- **Data flow edges**
- **Additional edges**
  - **Function Call -> Function Declaration**
  - **Variable Usage -> Variable Declaration**
  - **Field Access -> Field Declaration**
  - **Type -> Super Type**
  - **Expression/Variable -> Type**

## Why do we need extra pass?

1. Build AST file by file
2. Inter-file edges
3. Postprocess after building AST



# Introduction to Gremlin Query

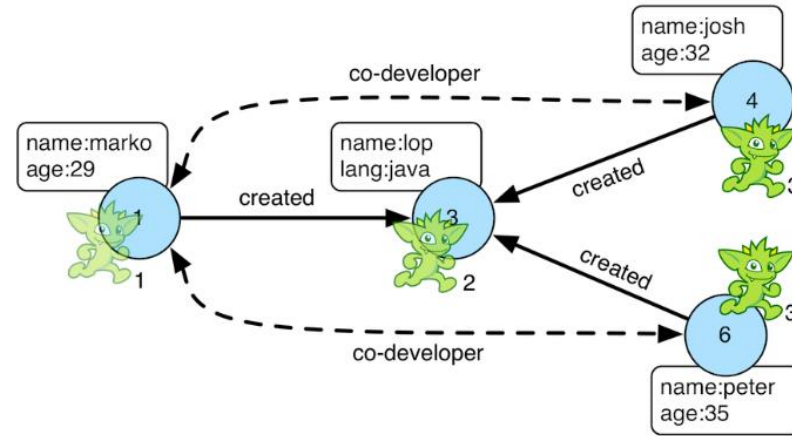
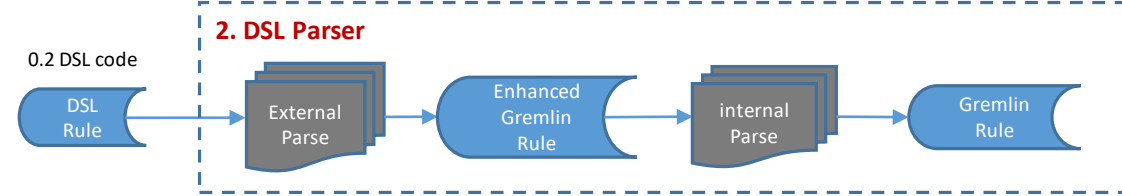
- Start steps: V() and E() are for query
  - addE(): Add edge to start the traversal
  - addV(): Add vertex to start the traversal
  - inject(): Inserts arbitrary objects to start the traversal
  - E(): Reads edges from the graph to start the traversal**
  - V(): Reads vertices from the graph to start the traversal**

## Five kinds of general step

- MapStep
  - FlatMapStep
  - FilterStep
  - SideEffectStep
  - BranchStep
- 

## End steps

- hasNext(): determines whether there are available results
- next(): returns the next result
- next(n): returns the next n results
- tryNext(): computes the next result
- toList(): return all results
- toSet(): return all results as a set
- toBulkSet(): return all results as a bulk set
- fill(collection): put results in provided collection
- iterate(): generate side effects without returning real result



## Query Example

1. Find the co-developers of marko

```
g.V().has("name", "marko").out("created").in("created")
    .has("name", P.neq("marko")).toList()
```

2. Find all developers older than 30

```
g.V().has("lang", "java").in("created")
    .where(__.has("age", P.gt(30))).toSet()
```



# DSL Translation

- Define our own DSL
  - Grammar of self-defined DSL

```
queryStmt: rootNodeAttr Satisfy condExpr Semicolon;

rootNodeAttr: (Node | NodeAttr) (Point NodeAttr)* (As alias | alias)?;

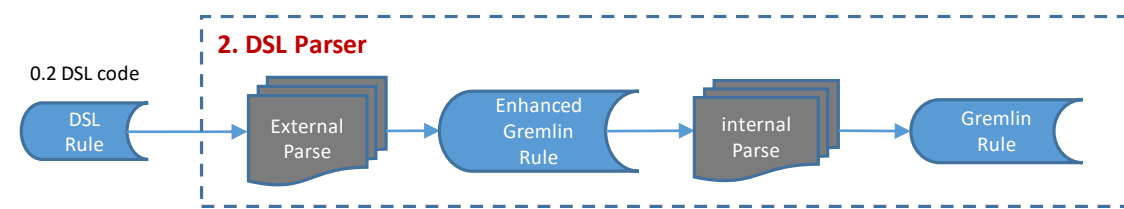
condExpr: (AND | OR) LeftBracket condExpr (Comma condExpr)+ RightBracket
          | NOT LeftBracket condExpr RightBracket
          | condition;
```

- Lexer of self-defined DSL

```
Node: 'functionCall' | 'functionDeclaration';

NodeAttr: 'function' | 'base' | 'body'
          | 'arguments' | 'rootBase'
          | 'enclosingClass';

Semicolon: ';;';
Point: '.';
```



- Enhanced Gremlin (APIs for JAVA)

```
@GremlinDsl(traversalSource = QueryTraversalSourceDsl")
public interface QueryTraversalDsl<S, E> extends GraphTraversal.Admin<S, E> {
    default GraphTraversal<S, Vertex> function() {
        return this.outE(Elabel.CALL_EDGE).inV();
    }

    default GraphTraversal<S, Vertex> functionCall() {
        return (GraphTraversal<S, Vertex>) this.hasLabel(P.within(LabelGroup.FUNCTION_CALL_LABEL));
    }

    default GraphTraversal<S, Vertex> functionDeclaration() {
        return (GraphTraversal<S, Vertex>) hasLabel(P.within(LabelGroup.FUNCTION_LABEL));
    }
}
```

- Original Gremlin



# Example of Multi-Layered DSL

- No debug code

- It is a Function
- Its name starts with "debug"
- It has only one parameter
- The type of parameter is java.util.List"

- Original Gremlin

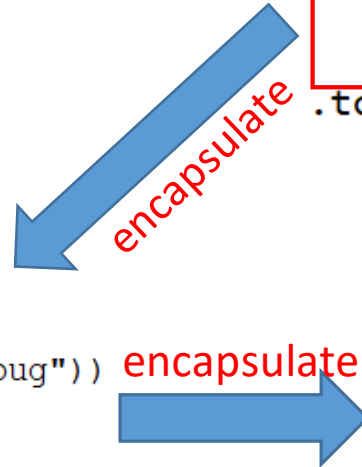
```
List<Vertex> vertexList = (List<Vertex>) this.graphTraversal
    .hasLabel(P.within(AbstractQuery.functionLabel))
    .has(VKey.NAME, TextP.startingWith("debug"))
    .where(__.outE(ELabel.AST_EDGE).has(EKey.RELATION,
        EValue.FUNC_DECL_PARAMETER)
        .inV().count().is(1))
    .where(__.outE(ELabel.AST_EDGE).has(EKey.RELATION,
        EValue.FUNC_DECL_PARAMETER)
        .inV().outE(ELabel.AST_EDGE).has(EKey.RELATION,
        EValue.NODE_TYPE)
        .inV().has("name", "java.util.List"))
    .toList();
```

- Enhanced Gremlin

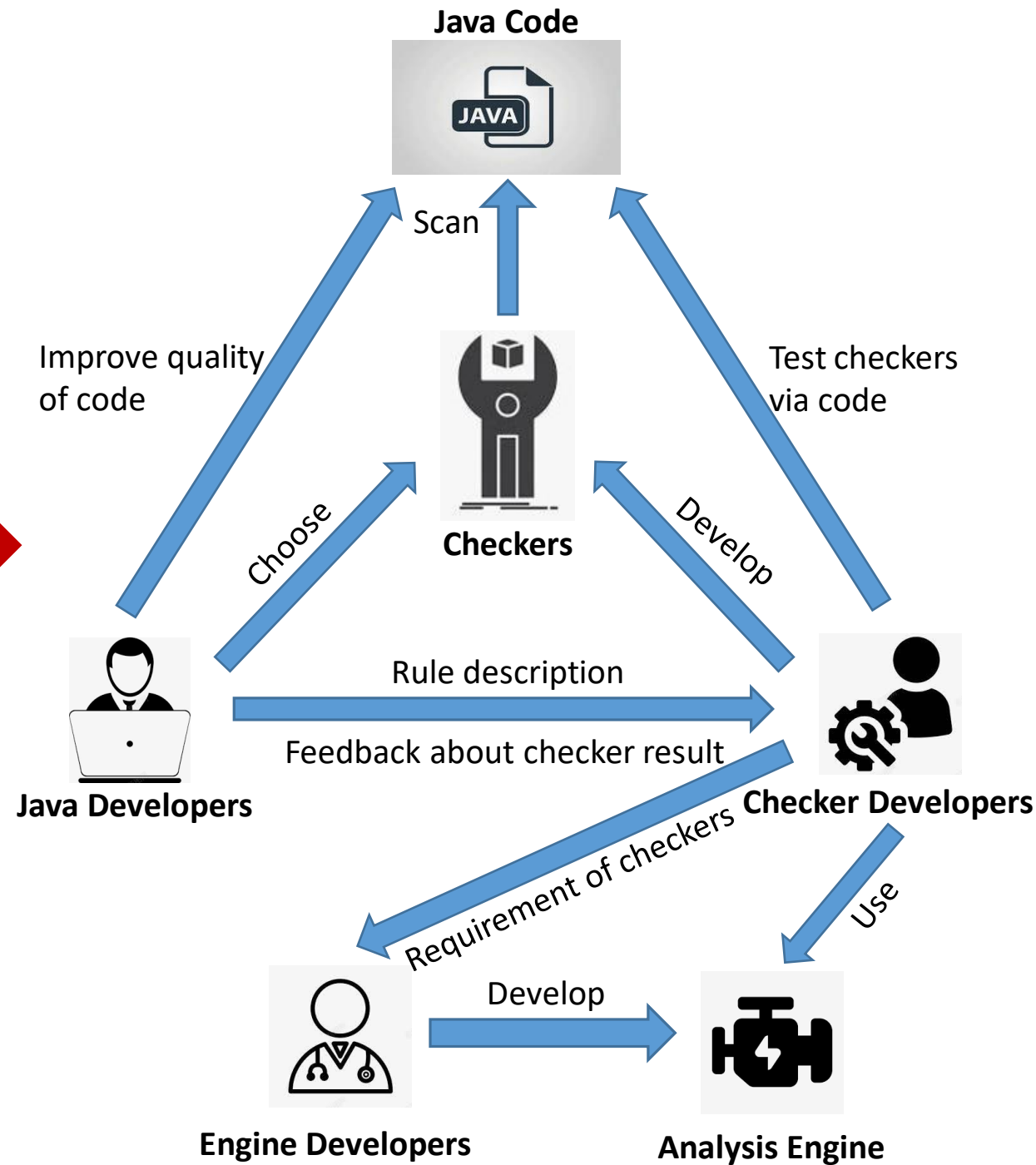
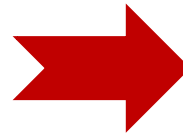
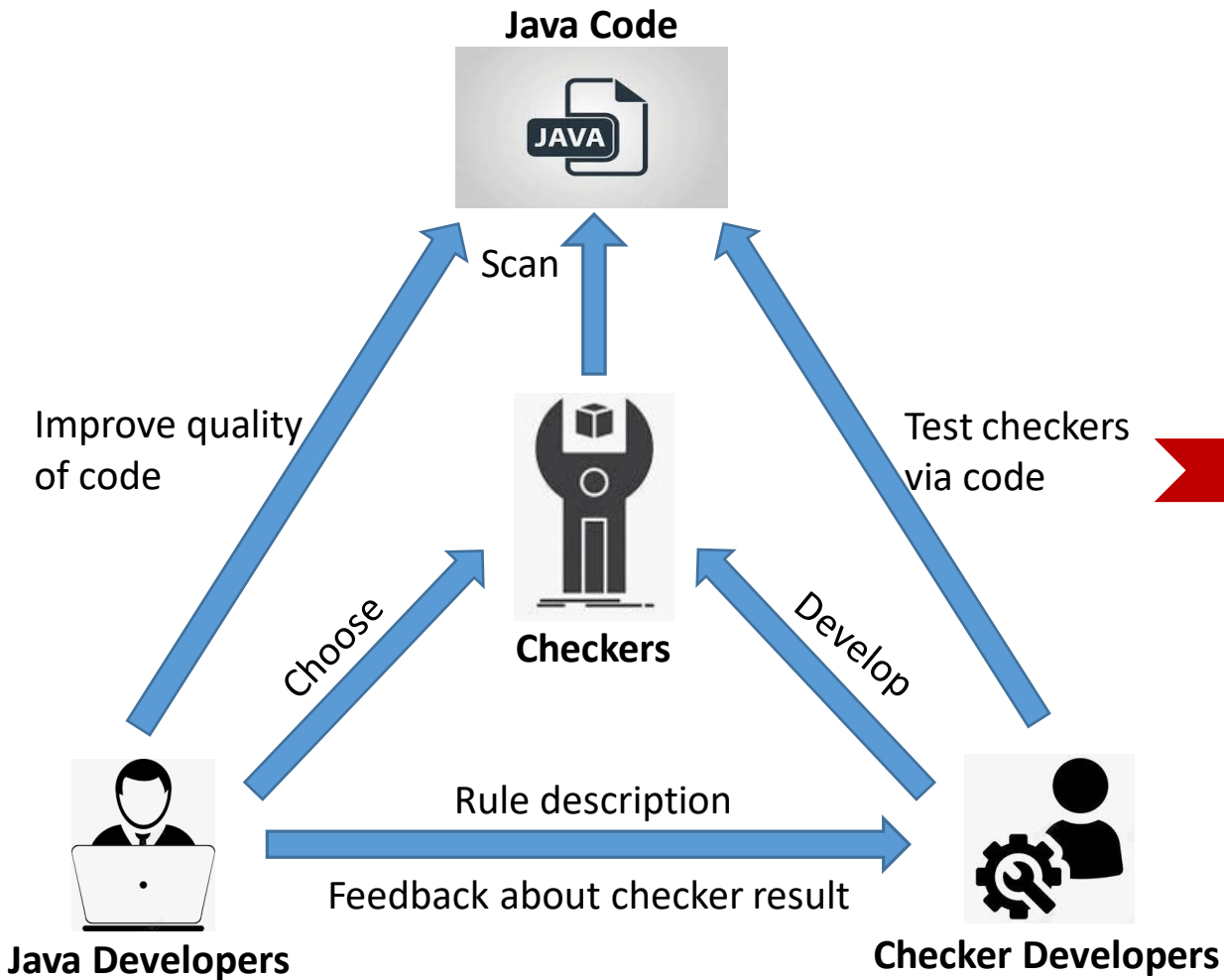
```
public List<Vertex> leftOverDebug() {
    return graphTraversal
        .function()
        .has("name", TextP.startingWith("debug"))
        .where(__.argument().count().is(1))
        .where(__.argument().type()
            .has("name", "java.util.List"))
        .toList();
}
```

- Self-defined DSL

```
functionDeclaration where
    and(name startWith "debug",
        arguments.size() == 1,
        arguments[0].type.name == "java.util.List");
```

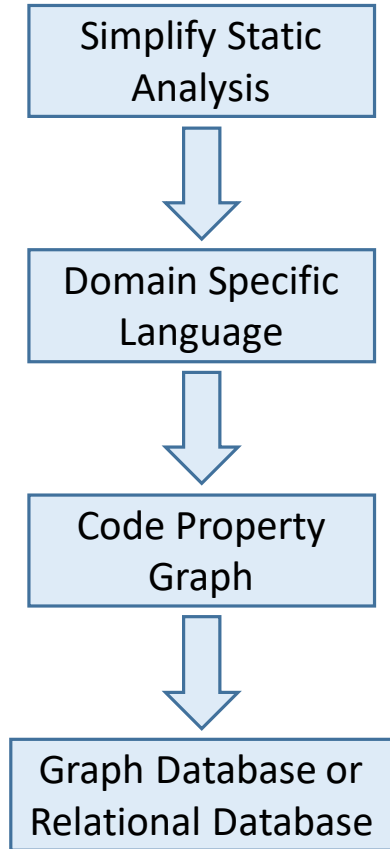


# Future of Static Analysis

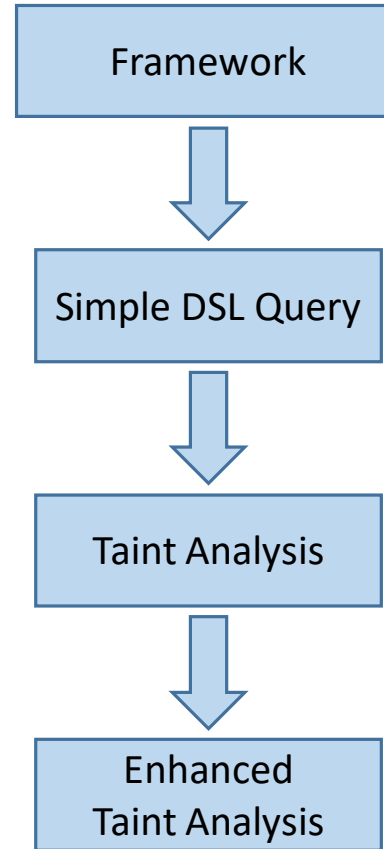


# Summary

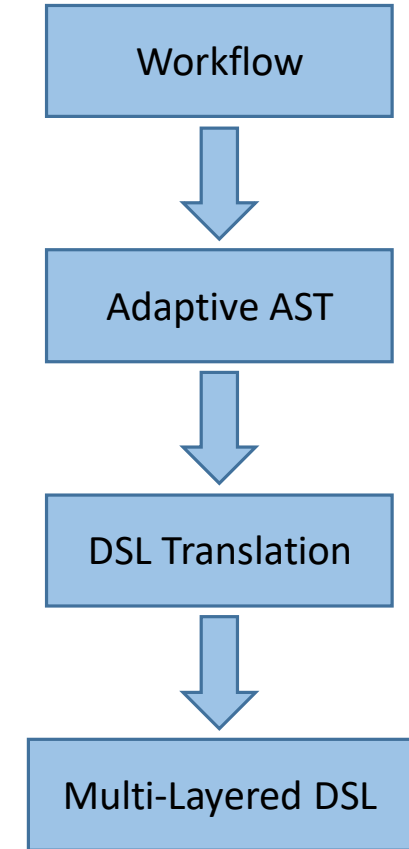
- Background



- CodeQL



- Our work



Thank You!

Спасибо!

谢谢！