

PyCG: Practical Call Graph Generation in Python

HE PEILIN



2021 IEEE/ACM 43rd International
Conference on Software Engineering
(ICSE)[1]

Contents

- 1 Introduction
- 2 Background
- 3 Topic
- 4 Evaluation and Implementation
- 5 References

Call Graph

Program Analysis, Inter-procedural Analysis, and Vulnerability Propagation Analysis.

Contributions

- *Static approach for pragmatic call graph generation in Python.*
- *Micro-benchmark suite used as a standard to evaluate this methods in Python.*
- *Evaluating the effectiveness of the approach through Micro-benchmark and Macro-benchmarks*
- *How the approach can aid dependency impact analysis through a potential enhancement of GitHub's "security advisory"*

Contents

- 1 Introduction
- 2 Background
- 3 Topic
- 4 Evaluation and Implementation
- 5 References

Contents

- 1 Introduction
- 2 Background
 - Challenges
 - Limitations of Existing Static Approaches
- 3 Topic
- 4 Evaluation and Implementation
- 5 References

Challenges List

Python Features

- Higher-order Functions
- Nested Definitions
- Classes
 - Inherit attributes and methods
 - Method Resolution Order (MRO)
- Modules
- Dynamic Features
 - Meta-programming
- Duck Typing

Contents

- 1 Introduction
- 2 Background
 - Challenges
 - Limitations of Existing Static Approaches
- 3 Topic
- 4 Evaluation and Implementation
- 5 References

crypto module

crypto.py

```
import cryptops

class Crypto:
    def __init__(self, key):
        self.key = key

    def apply(self, msg, func):
        return func(self.key, msg)

crp=Crypto('secretkey')
encrypted=crp.apply('hello_world', cryptops.encrypt)
decrypted=crp.apply(encrypted, cryptops.decrypt)
```

1

Call graphs for *crypto* module

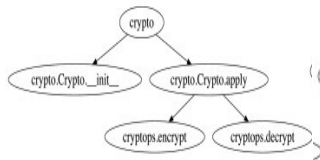


Figure 1. Precise call graph.

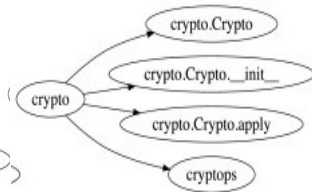


Figure 2. Pyan-generated call graph



Figure 3. Depends-generated call graph.

Contents

- 1 Introduction
- 2 Background
- 3 Topic
- 4 Evaluation and Implementation
- 5 References

Contents

- 1 Introduction
- 2 Background
- 3 Topic
 - The Core Analysis
 - Call Graph Construction
- 4 Evaluation and Implementation
- 5 References

Contents

1 Introduction

2 Background

3 Topic

■ The Core Analysis

■ Syntax

■ State

■ Analysis Rules

■ Call Graph Construction

4 Evaluation and Implementation

5 References

Syntax Analysis – AST

$$e \in Expr ::= o \mid x \mid x := e \mid \mathbf{function} \ x \ (y \dots) \ e \mid \mathbf{return} \ e \mid \\ e(x=e \dots) \mid \mathbf{class} \ x \ (y \dots) \ e \mid e.x \mid e.x := e \mid \\ \mathbf{new} \ x \ (y = e \dots) \mid \mathbf{import} \ x \ \mathbf{from} \ m \ \mathbf{as} \ y \mid \\ \mathbf{iter} \ x \mid e;e$$
$$o \in Obj ::= n, v$$
$$v \in Definition ::= x, \tau$$
$$\tau \in IdentType ::= \mathbf{func} \mid \mathbf{var} \mid \mathbf{cls} \mid \mathbf{mod}$$
$$n \in Namespace ::= (v)^*$$
$$x, y \in Identifier ::= \text{is the set of program identifiers}$$
$$m \in Modules ::= \text{is the set of modules}$$

Figure 4. The syntax for representing the input Python programs along with the evaluation contexts

Evaluation E

Use evaluation contexts [3], [4] that describe the order in which sub-expressions are evaluated.

$$\begin{aligned} E ::= & [] \mid x := E \mid \text{return } E \mid E(x = e \dots) \mid \\ & o(x = E \dots) \mid \text{new } x(y = E) \mid E.x \mid E.x := e \mid \\ & o.x := E \mid \text{iter } o \mid E; e \mid o; E \end{aligned}$$

Contents

1 Introduction

2 Background

3 Topic

■ The Core Analysis

■ Syntax

■ State

■ Analysis Rules

■ Call Graph Construction

4 Evaluation and Implementation

5 References

After the AST (Figure.4), the analysis maintains a state consisting of four domains as shown:

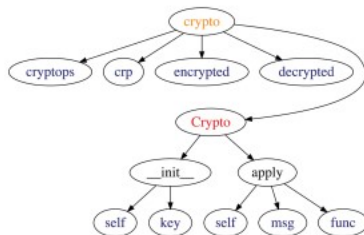
$$\pi \in \text{AssignG} = \text{Obj} \hookrightarrow P(\text{Obj})$$

$$s \in \text{Scope} = \text{Definition} \hookrightarrow P(\text{Definition})$$

$$h \in \text{ClassHier} = \text{Obj} \hookrightarrow \text{Obj}^*$$

$$\sigma \in \text{State} = \text{AssignG} \times \text{Scope} \times \text{Namespace} \times \text{ClassHier}$$

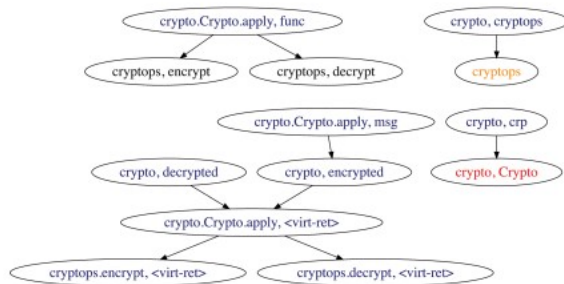
Scope Tree



(a) The scope tree of the `crypto` module.

Figure 5. The scope tree of the `crypto` module

Assignment Graph



(b) The assignment graph of the `crypto` module.

Figure 6. The assignment graph of the `crypto` module

Contents

1 Introduction

2 Background

3 Topic

■ The Core Analysis

- Syntax

- State

- Analysis Rules

- Call Graph Construction

4 Evaluation and Implementation

5 References

Analysis Rules

Demonstrating the state transition rules of analysis, the rules follow the form:

$$\langle \pi, s, n, h, E[e] \rangle \rightarrow \langle \pi', s', n', h', E[e'] \rangle \quad (1)$$

When having an expression e in the evaluation context E

$$\text{E-CTX} : \frac{\langle \pi, s, n, h, E[e] \rangle \hookrightarrow \langle \pi', s', n', h', e' \rangle}{\langle \pi, s, n, h, E[e] \rangle \rightarrow \langle \pi', s', n', h', E[e'] \rangle} \quad (2)$$

getObject(s, n, x)

The function iterates every element y of the namespace n in the reverse order, checking whether the element node y has any child matching the identifier x

addScope(s, n, x, t)

This function adds an edge from the node accessed by the path n to the target node given by the definition (x, t)

getClassAttrObject(o, x, h)

This function deals with multiple inheritance, retrieving the object corresponding to the attribute x of the receiver object o , the analysis by examining the hierarchy of classes h

Python Source Code

main.py

```
class A:
    def func():
        pass

class B:
    def func():
        pass

a=A()
b=B()
a.func()
b.func()
```

main.py

```
1 class A:
2     def func():
3         pass
4
5 class B:
6     def func():
7         pass
8
9 class C(B,A):
10     pass
11
12 c=C()
13 c.func()
```

COMPOUND

$$\frac{}{\langle \pi, s, n, h, E[o_1; o_2] \rangle \rightarrow \langle \pi, s, n, h, E[o_2] \rangle}$$

IDENT

$$\frac{o = \text{getObject}(s, n, x)}{\langle \pi, s, n, h, E[x] \rangle \rightarrow \langle \pi, s, n, h, E[o] \rangle}$$

ASSIGN

$$\frac{s' = \text{addScope}(s, n, x, \text{var}) \quad o' = (n, (x, \text{var})) \quad \pi' = \pi[o' \rightarrow \pi(o') \cup \{o\}]}{\langle \pi, s, n, h, E[x := o] \rangle \rightarrow \langle \pi', s', n, h, E[o'] \rangle}$$

FUNC

$$\frac{s' = \text{addScope}(s, n, x, \text{func}) \quad n' = n \cdot (x, \text{func}) \quad s'' = \text{addScope}(s', n', \text{ret}, \text{var}) \quad s^{(3)} = \text{addScope}(s'', n', y, \text{var})}{\langle \pi, s, n, h, E[\text{function } x (y \dots) e] \rangle \rightarrow \langle \pi, s^{(3)}, n', h, E[e] \rangle}$$

RETURN

$$\frac{o' = (n \cdot x, (\text{ret}, \text{var})) \quad \pi' = \pi[o' \rightarrow \pi(o') \cup \{o\}]}{\langle \pi, s, n \cdot x, h, E[\text{return } o] \rangle \rightarrow \langle \pi', s, n, h, E[o'] \rangle}$$

CALL

$$\frac{o_1 = (n', (f, \text{func})) \quad o'_2 = (n' \cdot f, (y, \text{var})) \quad \pi' = \pi[o'_2 \rightarrow \pi(o'_2) \cup \{o_2\}]}{\langle \pi, s, n, h, E[o_1(y = o_2 \dots)] \rangle \rightarrow \langle \pi', s, n, h, (n' \cdot f, (\text{ret}, \text{var})) \rangle}$$

Figure 7. Rules of the analysis (1)

Analysis Rules

$$\begin{array}{c} \text{CLASS} \\ \frac{s' = \text{addScope}(s, n, x, \text{cls}) \quad t = \langle \text{getObject}(s, n, b) \mid b \in (y \dots) \rangle \\ \quad h' = h[(n, (x, \text{cls})) \rightarrow t] \quad n' = n \cdot (x, \text{cls})}{\langle \pi, s, n, h, E[\text{class } x (y \dots) e] \rangle \rightarrow \langle \pi, s', n', h', E[e] \rangle} \\ \\ \text{ATTR} \\ \frac{o' = \text{getClassAttrObject}(o, x, h)}{\langle \pi, s, n, h, E[o.x] \rangle \rightarrow \langle \pi, s, c, h, E[o'] \rangle} \\ \\ \text{NEW} \\ \frac{o_3 = \text{getObject}(s, n, x) \\ o_2 = \text{getClassAttrObject}(o_3, __\text{init}__, h)}{\langle \pi, s, n, h, E[\text{new } x(y = o_1 \dots)] \rangle \rightarrow \langle \pi, s, n, h, E[o_2(y = o_1 \dots); o_3] \rangle} \\ \\ \text{ATTR-ASSIGN} \\ \frac{o_3 = \text{getClassAttrObject}(o_1, x, h) \quad \pi' = \pi[o_3 \rightarrow \pi(o_3) \cup \{o_2\}]}{\langle \pi, s, n, h, E[o_1.x := o_2] \rangle \rightarrow \langle \pi', s, n, h, E[o_3] \rangle} \\ \\ \text{IMPORT} \\ \frac{o_2 = \text{getObject}(s, m, x) \quad s' = \text{addScope}(s, n, y, \text{var}) \\ o_1 = (n, (y, \text{var})) \quad \pi' = \pi[o_1 \rightarrow \pi(o_1) \cup \{o_2\}]}{\langle \pi, s, n, h, E[\text{import } x \text{ from } m \text{ as } y] \rangle \rightarrow \langle \pi', s', n, h, E[o_1] \rangle} \\ \\ \text{ITER-ITERABLE} \\ \frac{o' = \text{getClassAttrObject}(o, __\text{next}__, h)}{\langle \pi, s, n, h, E[\text{iter } o] \rangle \rightarrow \langle \pi, s, n, h, E[o'()] \rangle} \\ \\ \text{ITER-GENERATOR} \\ \frac{\text{getClassAttrObject}(o, __\text{next}__, h) = \text{undefined}}{\langle \pi, s, n, h, E[\text{iter } o] \rangle \rightarrow \langle \pi, s, n, h, E[o()] \rangle} \end{array}$$

Figure 8. Rules of the analysis (2)

Contents

- 1 Introduction
- 2 Background
- 3 Topic
 - The Core Analysis
 - Call Graph Construction
- 4 Evaluation and Implementation
- 5 References

Producing a call graph:

$$cg \in CallGraph = Obj \hookrightarrow P(Obj) \quad (3)$$

Algorithm (Figure.9) takes two elements as input:

- 1 a program $p \in Program$ of the model language whose syntax is shown in Figure.4.
- 2 the final state $\sigma \in State$ stemming from the analysis step.

Algorithm

```
Input  :  $p \in \text{Program}$   
          $\sigma \in \text{State}$   
Output:  $cg \in \text{CallGraph}$   
1 foreach  $e$  in  $\text{Program}$  do  
2   while  $e \notin \text{Obj}$  do  
3      $\langle \sigma, E[e] \rangle \rightarrow \langle \sigma', E[e'] \rangle$   
4     if  $e' = o_1(y = o_2 \dots)$  then // Call Expression  
5        $(\pi, s, n \cdot f, h) \leftarrow \sigma'$   
6        $c \leftarrow \text{getReachableFuns}(\pi, o_1)$   
7        $o_3 \leftarrow \text{getObject}(s, n, f)$   
8        $cg \leftarrow cg[o_3 \rightarrow cg(o_3) \cup c]$  // Add Call Edges  
9     end  
10     $e \leftarrow e'$   
11  end  
12 end  
13 return  $cg$ 
```

Figure 9. Algorithm for call graph construction.

Contents

- 1 Introduction
- 2 Background
- 3 Topic
- 4 Evaluation and Implementation
- 5 References

Evaluating the approach based on three research questions:

- RQ1 Is the proposed approach effective in constructing call graphs for Python programs?
- RQ2 How does the proposed approach stand in comparison with existing open-source, static-based approaches for Python?
- RQ3 What is the performance of our approach?

Setup

Experiments on a **Debian 9 host** with 16 CPUs and 16 GBs of RAM.

- 1 a micro-benchmark suite containing 112 minimal Python programs.
- 2 a macro-benchmark suite of five popular real-world Python packages.

Contents

- 1 Introduction
- 2 Background
- 3 Topic
- 4 Evaluation and Implementation
 - Micro-benchmark
 - Macro-benchmark
 - Time and Memory Performance
 - Case Study
- 5 References

Micro-benchmark

Category	#tests	Description
parameters	6	Positional arguments that are functions
assignments	4	Assignment of functions to variables
built-ins	3	Calls to built in functions and data types
classes	22	Class construction, attributes, methods
decorators	7	Function decorators
dicts	12	Hashmap with values that are functions
direct calls	4	Direct call of a returned function (<code>func()()</code>)
exceptions	3	Exceptions
functions	4	Vanilla function calls
generators	6	Generators
imports	14	Imported modules, functions classes
kwargs	3	Keyword arguments that are functions
lambdas	5	Lambdas
lists	8	Lists with values that are functions
mro	7	Method Resolution Order (MRO)
returns	4	Returns that are functions

Figure 10. Micro-benchmark for the suite categories.

Addressing Validity Threats

Asking two Python developers to rank the suite (from 1 to 10) based on the following criteria:

- 1 Completeness : Does it cover Python features?
- 2 Code Quality : Are the tests unique and minimal?
- 3 Description : Does the description adequately describe the given test case?

Micro-benchmark Suite Results

Category	PyCG		Pyan	
	Complete	Sound	Complete	Sound
assignments	4/4	3/4	4/4	4/4
built-ins	3/3	1/3	2/3	0/3
classes	22/22	22/22	6/22	10/22
decorators	6/7	5/7	4/7	3/7
dicts	12/12	11/12	6/12	6/12
direct calls	4/4	4/4	0/4	0/4
exceptions	3/3	3/3	0/3	0/3
functions	4/4	4/4	4/4	3/4
generators	6/6	6/6	0/6	0/6
imports	14/14	14/14	10/14	4/14
kwargs	3/3	3/3	0/3	0/3
lambdas	5/5	5/5	4/5	0/5
lists	8/8	7/8	3/8	4/8
mro	7/7	5/7	0/7	2/7
parameters	6/6	6/6	0/6	0/6
returns	4/4	4/4	0/4	0/4
Total	111/112	103/112	43/112	36/112

Figure 11. Micro-benchmark results for PyCG and Pyan. Depends is unsound in all cases and complete in 110/112 cases and is omitted.

Contents

- 1 Introduction
- 2 Background
- 3 Topic
- 4 Evaluation and Implementation**
 - Micro-benchmark
 - **Macro-benchmark**
 - Time and Memory Performance
 - Case Study
- 5 References

Macro-benchmark

Project	LoC	Stars	Forks	Description
fabric	3,236	12.1k	1.8k	Remote execution & deployment
autojump	2,662	10.8k	530	Directory navigation tool
asciinema	1,409	7.9k	687	Terminal session recorder
face_classification	1,455	4.7k	1.4k	Face detection & classification
Sublist3r	1,269	4.4k	1.1k	Subdomains enumeration tool

Figure 12. Macro-benchmark suite project details.

Macro-benchmark Results

Project	Precision (%)			Recall (%)		
	PyCG	Pyan	Depends	PyCG	Pyan	Depends
autojump	99.5	66.5	99.2	68.2	28.5	22.5
fabric	98.3	-	100	61.9	-	6.3
asciinema	100	-	98.1	68	-	15.5
face_classification	99.5	86.8	96.2	89.7	7.6	5.7
Sublist3r	98.8	69.8	100	61.6	25.6	21.9
Average	99.2	74.4	98.7	69.9	20.6	14.4

Figure 13. Macro-benchmark results and tool comparison.

Contents

- 1 Introduction
- 2 Background
- 3 Topic
- 4 Evaluation and Implementation**
 - Micro-benchmark
 - Macro-benchmark
 - **Time and Memory Performance**
 - Case Study
- 5 References

Time and Memory Performance

Command

UNIX command : *time, pmap* (average out of 20 runs)

Project	Time (sec)			Memory (MB)		
	PyCG	Pyan	Depends	PyCG	Pyan	Depends
autojump	0.76	0.42	2.37	62.7	37.8	27.1
fabric	0.77	-	1.83	60.9	-	18.5
asciinema	0.87	-	2	61.6	-	19.4
face_classification	0.92	0.38	2.49	60.9	35.3	25.6
Sublist3r	0.51	0.33	2.01	60	35.8	19.4
Average	0.77	0.38	2.14	61.2	36.3	22

Figure 14. Time and Memory Performance.

Contents

- 1 Introduction
- 2 Background
- 3 Topic
- 4 Evaluation and Implementation**
 - Micro-benchmark
 - Macro-benchmark
 - Time and Memory Performance
 - **Case Study**
- 5 References

Case Study: A Fine-grained Tracking of Vulnerable Dependencies

Further, we show a potential application through the enhancement of GitHub's "security advisory"[5] notification service.

Cases contained functional vulnerability

- 1 **PyYAML**[6] (versions before 5.1), a YAML parser affected by CVE-2017-18342 [7].
- 2 **Paramiko**[8] (multiple versions before 2.4.1), an implementation of the SSHv2 protocol affected by CVE-2018-7750 [9].

Results

- The vulnerable function in PyYAML(i.e., *load*) was invoked by 42/106 projects.
- In Paramiko, method (*start_server*) was not utilized at all by any of the 76 projects. We also observed that 12 projects did not invoke any library coming from Paramiko.

Contents

- 1 Introduction
- 2 Background
- 3 Topic
- 4 Evaluation and Implementation
- 5 References

References (1)

Paper

- [1] Vitalis Salis et al. “PyCG: Practical Call Graph Generation in Python”. In: *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. 2021, pp. 1646–1657. DOI: 10.1109/ICSE43902.2021.00146.

Related Work and Additional Literature

- [3] Matthias Felleisen, Robert Bruce Findler, and Matthew Flatt. *Semantics engineering with PLT Redex*. Mit Press, 2009.
- [4] Magnus Madsen, Ondřej Lhoták, and Frank Tip. “A model for reasoning about JavaScript promises”. In: *Proceedings of the ACM on Programming Languages* 1.OOPSLA (2017), pp. 1–24.

References (2)

Website Sources

Title page: drocheam. *About LaTeX Beamer Template in TH Koeln Style*. Online.
<https://github.com/drocheam/th-koeln-beamer-template>

- [5] *GitHub advisory database*. Online. <https://github.com/advisories>. 2020.
- [6] *PyYAML: The next generation YAML parser and emitter for Python*. Online.
<https://github.com/yaml/pyyaml/>. 2020.
- [7] *CVE-2017-18342*. Online. <https://nvd.nist.gov/vuln/detail/CVE-2017-18342>. 2017.
- [8] *Paramiko: The leading native Python SSHv2 protocol library*. Online.
<https://github.com/paramiko/paramiko/>. 2020.
- [9] *CVE-2018-7750*. Online. <https://nvd.nist.gov/vuln/detail/CVE-2018-7750>. 2018.

Thank you for your attention

Please feel free to ask any questions

Contact:

`timhh991022@gmail.com`