



THE UNIVERSITY OF TEXAS AT DALLAS

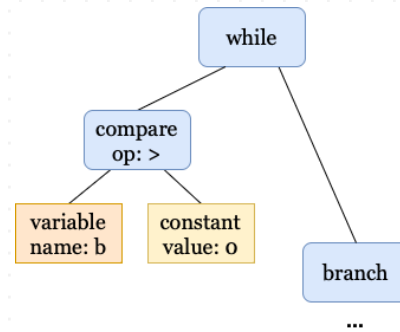
(Partial) Program Dependence Learning

Aashish Yadavally and Tien N. Nguyen
Computer Science Department
The University of Texas at Dallas

Wenbo Wang and Shaohua Wang
Department of Informatics
New Jersey Institute of Technology

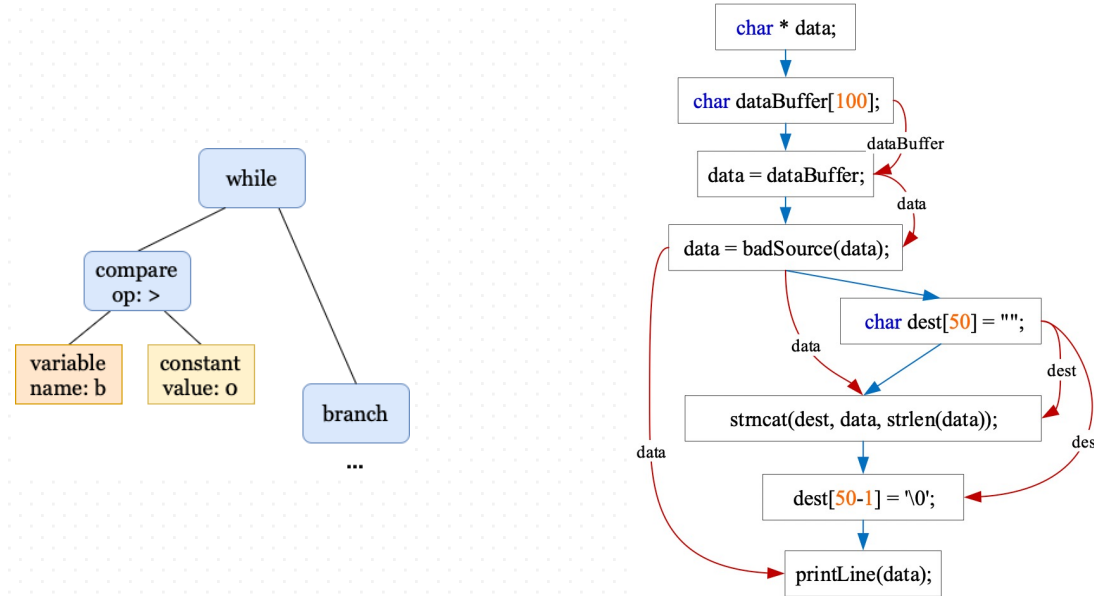
Background

- ❖ Given complete program units, one can utilize different tools to help build program representations such as abstract syntax trees (ASTs)



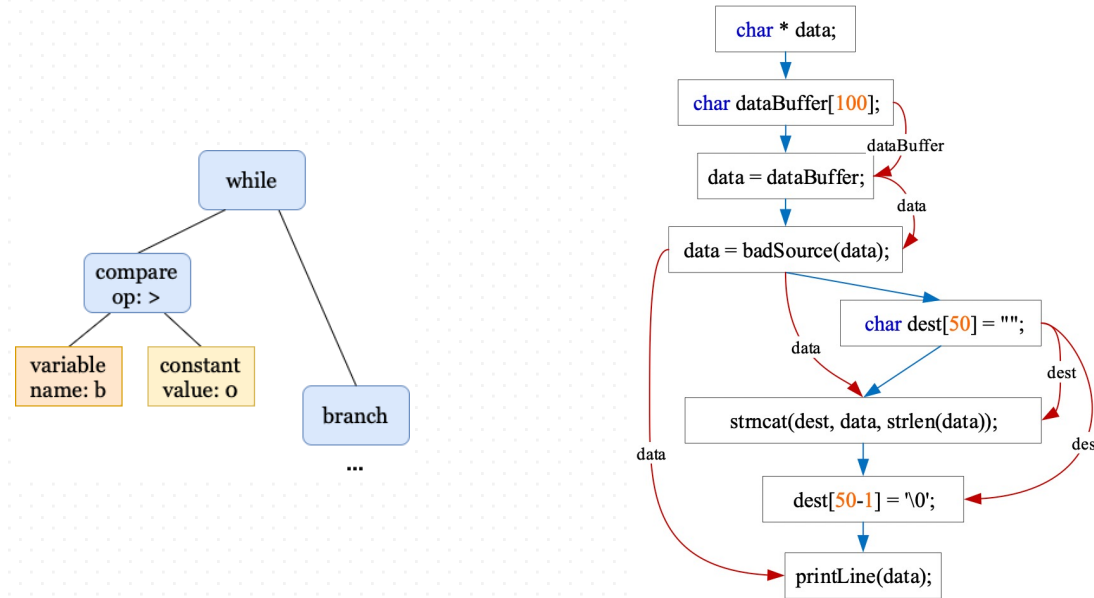
Background

- ❖ Given complete program units, one can utilize different tools to help build program representations such as abstract syntax trees (ASTs), control-flow and program dependence graphs (CFGs & PDGs), etc.



Background

- ❖ Given complete program units, one can utilize different tools to help build program representations such as abstract syntax trees (ASTs), control-flow and program dependence graphs (CFGs & PDGs), etc.



- ❖ In contrast, analyzing code fragments from online forums is difficult as they are often incomplete, unparseable, contain declaration/reference ambiguity, and are interspersed between user comments.

Let us call this *partial program dependence analysis*.

Why is this useful?

Why is this useful?

An Empirical Study of C++ Vulnerabilities in Crowd-Sourced Code Examples

Morteza Verdi, Ashkan Sami, Jafar Akhondali, Foutse Khomh, Gias Uddin, and Alireza Karami Motlagh

Abstract—Software developers share programming solutions in Q&A sites like Stack Overflow, Stack Exchange, Android forum, and so on. The reuse of crowd-sourced code snippets can facilitate rapid prototyping. However, recent research shows that the shared code snippets may be of low quality and can even contain vulnerabilities. This paper aims to understand the nature and the prevalence of security vulnerabilities in crowd-sourced code examples. To achieve this goal, we investigate security vulnerabilities in the C++ code snippets shared on Stack Overflow over a period of 10 years. In collaborative sessions involving multiple human coders, we manually assessed each code snippet for security vulnerabilities following CWE (Common Weakness Enumeration) guidelines. From the 72,483 reviewed code snippets used in at least one project hosted on GitHub, we found a total of 99 vulnerable code snippets categorized into 31 types. Many of the investigated code snippets are still not corrected on Stack Overflow. The 99 vulnerable code snippets found in Stack Overflow were reused in a total of 2859 GitHub projects. To help improve the quality of code snippets shared on Stack Overflow, we developed a browser extension that allows Stack Overflow users to be notified for vulnerabilities in code snippets when they see them on the platform.

Index Terms—Stack Overflow, Software Security, C++, SOTorrent, Vulnerability Migration, GitHub, Vulnerability Evolution

Why is this useful?

- ❖ Some existing approaches build CFG/PDGs for partial code via *manual intervention* (e.g., by wrapping around method signature) in a best-effort manner, often at the cost of several misses.
-

Why is this useful?

- ❖ Some existing approaches build CFG/PDGs for partial code via *manual intervention* (e.g., by wrapping around method signature) in a best-effort manner, often at the cost of several misses.

```
1 string subTag(string s, string a, string b) {  
2   std::string lower_s;  
3   std::transform(s.begin(), s.end(), lower_s.begin(), ::  
4     tolower);  
5   std::transform(a.begin(), a.end(), a.begin(), ::tolower);  
6   auto position = lower_s.find(a);  
   ...  
}
```

In the above code listing, due to the unknown data type `string` on line 2 and API element `transform` on lines 3 and 4, traditional program analysis tools ignore all CFG/PDG edges to/from these statements. However, this listing is vulnerable on line 3, automated detection of which is not possible due to the missing edges.

Motivating Examples

Observation 1

Partial program dependence analysis is desirable for the software engineering tasks in which completely analyzable code is not available.

```
1 std :: shared_ptr<FILE> pipe(popen(cmd, "r"), pclose);
2 if (! pipe) return "ERROR";
3 char buffer[128];
4 std :: string result = "";
5 while (! feof (pipe.get() ) ) {
6     if ( fgets (buffer , 128, pipe.get() ) != NULL)
7         result += buffer;
8 }
```

Illustration 1: Incomplete S/O code snippet to execute a command within a C++ program, that is prone to OS Command Injection.

Motivating Examples

Observation 1

Partial program dependence analysis is desirable for the software engineering tasks in which completely analyzable code is not available.

```
1 #include <cstdio>
2 #include <iostream>
3 #include <memory>
4 #include <stdexcept>
5 #include <string>
6 #include <array>
7
8 std::string exec(const char* cmd) {
9     std::array<char, 128> buffer;
10    std::string result;
11    std::unique_ptr<FILE, decltype(&pclose)> pipe(popen(cmd,
12                                                    "r"), pclose);
13    if (!pipe) {
14        throw std::runtime_error("popen() failed!");
15    }
16    while (fgets(buffer.data(), buffer.size(),
17                pipe.get()) != nullptr) {
18        result += buffer.data();
19    }
20    return result;
21 }
```

Illustration 2: Complete code example with the same POSIX elements as in Illustration 1. Here, similar statement blocks are highlighted with same color.

```
1 std::shared_ptr<FILE> pipe(popen(cmd, "r"), pclose);
2 if (! pipe) return "ERROR";
3 char buffer[128];
4 std::string result = "";
5 while (! feof (pipe.get() ) ) {
6     if ( fgets (buffer , 128, pipe.get() ) != NULL)
7         result += buffer;
8 }
```

Illustration 1: Incomplete S/O code snippet to execute a command within a C++ program, that is prone to OS Command Injection.

Observation 2

Inter-statement dependence analysis of partial code can be derived from the patterns learned from such analyses of entire programs in existing code corpora.

Key Ideas

❖ PDGs are *repetitive*!

Nguyen *et al.* [1] reported that among 17.5M PDGs with 1.6B PDG subgraphs, 14.3% of the PDGs have all of their subgraphs repeated across different projects. Furthermore, in 15.6% of the PDGs, at least 90% of their subgraphs are likely to have appeared before in other projects.

Key Ideas

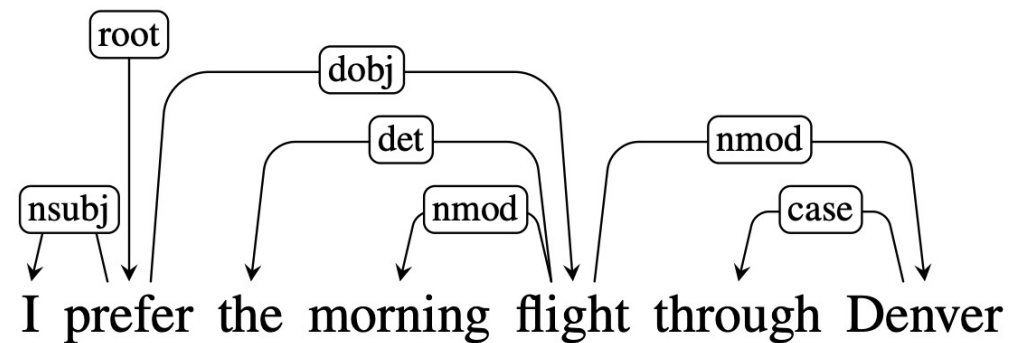
- ❖ PDGs are *repetitive*!
- ❖ We can leverage a pattern learning-based approach to **partial program dependence analysis**.

Key Ideas

- ❖ PDGs are *repetitive*!
- ❖ We can leverage a pattern learning-based approach to **partial program dependence analysis**.
- ❖ Is this even possible?

Key Ideas

- ❖ PDGs are *repetitive*!
- ❖ We can leverage a pattern learning-based approach to **partial program dependence analysis**.
- ❖ Is this even possible?



General Overview

```
1 private boolean isValidUntil (Until annotation) {  
2     if (annotation != null) {  
3         double annotationVersion = annotation.value();  
4         if (annotationVersion <= version) {  
5             return false;  
6         }  
7     }  
8     return true;  
9 }
```

Code Snippet

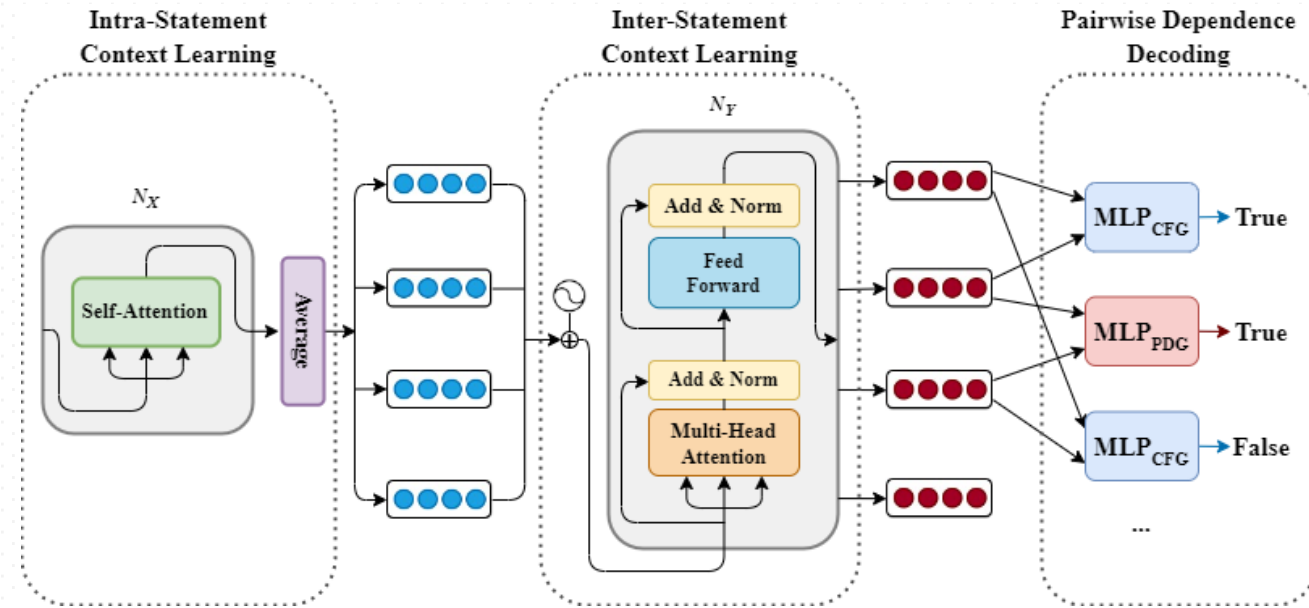
Training Process

- ❖ From *complete* code corpora, build CFG/PDGs to utilize as ground-truth.
- ❖ Train a model to predict the presence of a CFG/PDG edge between two statements in a specific code context.

Inference

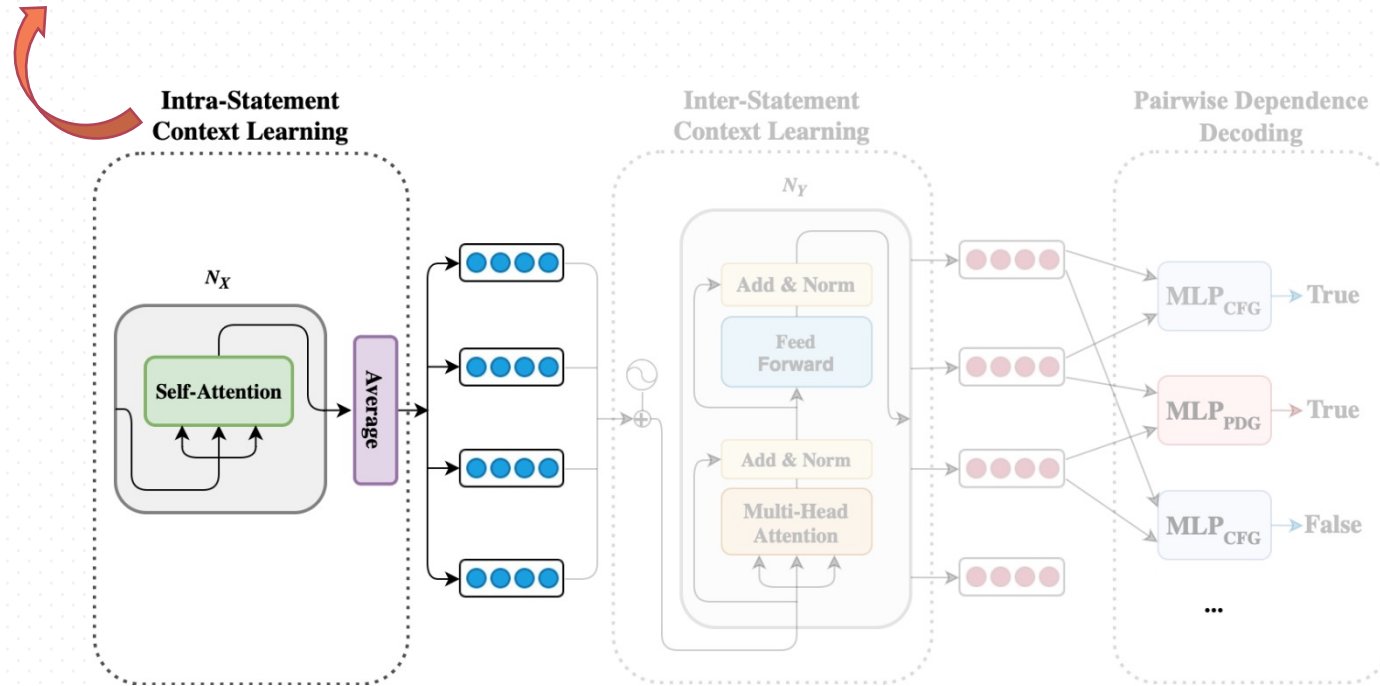
- ❖ Given *partial/complete* code snippet, predict the presence of a CFG/PDG edge between each statement pair.
- ❖ Combination of all such CFG/PDG edges can be realized as a CFG/PDG of the code snippet.

NEURALPDA: Neural Network-Based Program Dependence Analysis

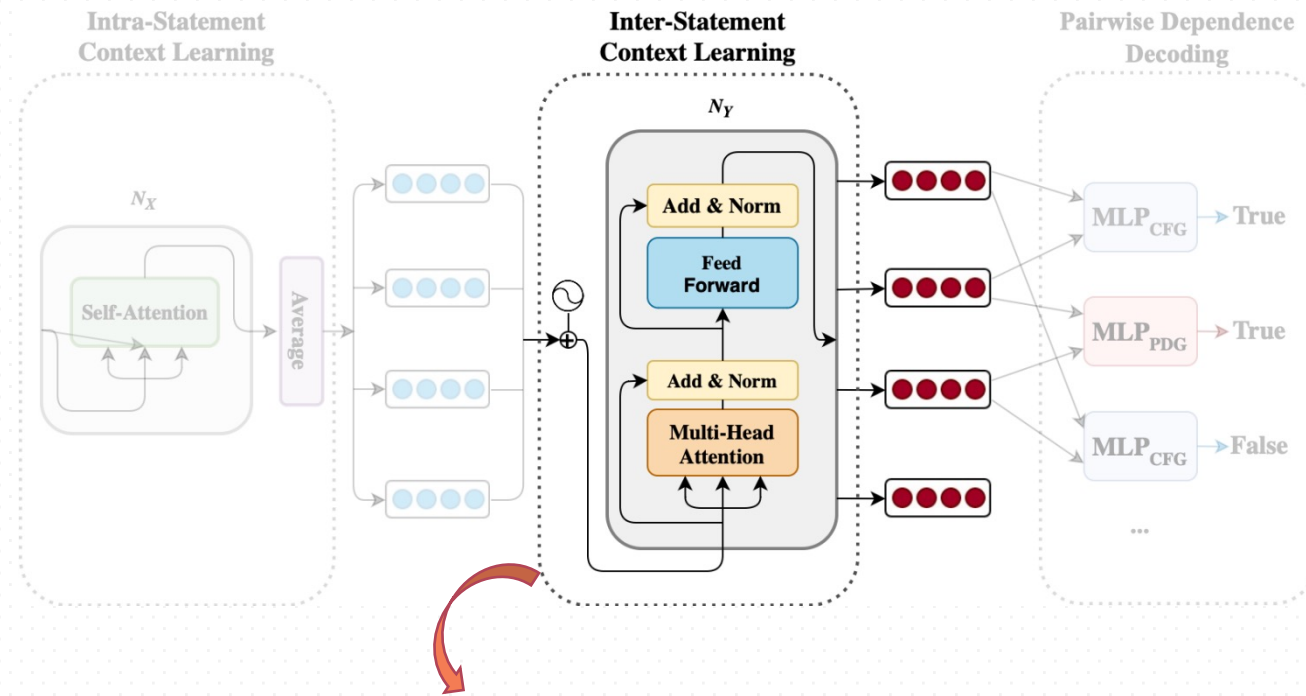


NEURALPDA: Neural Network-Based Program Dependence Analysis

Helps *relay the syntactic and semantic relationships* between the tokens in a statement to other statements in the code snippet.



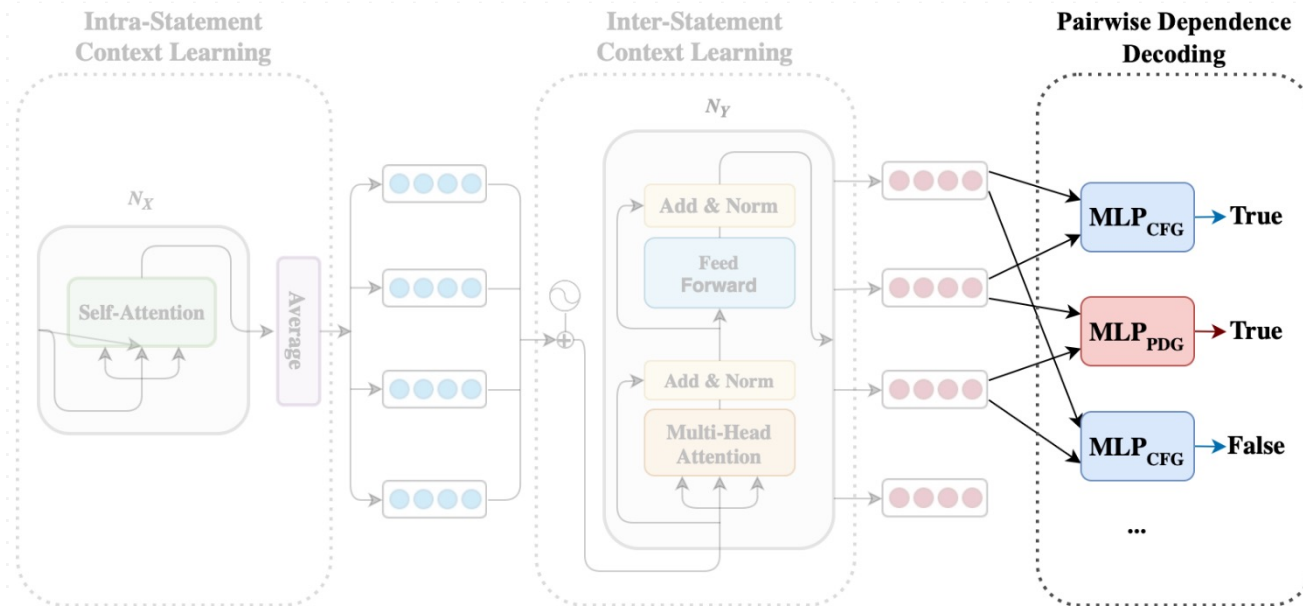
NEURALPDA: Neural Network-Based Program Dependence Analysis



The goal of this component is to *learn latent representations* for each statement that model the inter-statement dependencies.

NEURALPDA: Neural Network-Based Program Dependence Analysis

The combination of all the CFG/PDG edges extracted via such an arc-factored approach is *realized as the CFG/PDG* for the given program.



Empirical Evaluation (Intrinsic)

P/L	Graph	Accuracy	Precision	Recall	F1-Score
Java	<i>CFG</i>	99.79	98.31	98.58	98.44
	<i>PDG</i>	98.87	89.89	87.53	88.70
	<i>Overall</i>	99.33	94.75	93.83	94.29
C/C++	<i>CFG</i>	99.50	96.76	96.56	96.66
	<i>PDG</i>	98.55	83.55	90.01	86.66
	<i>Overall</i>	99.02	91.10	93.87	92.46

Table 1: Performance of NEURALPDA on **complete** code from Java and C/C++ (*Intrinsic Evaluation*)

Empirical Evaluation (Intrinsic)

P/L	Graph	Accuracy	Precision	Recall	F1-Score
Java	<i>CFG</i>	99.79	98.31	98.58	98.44
	<i>PDG</i>	98.87	89.89	87.53	88.70
	<i>Overall</i>	99.33	94.75	93.83	94.29
C/C++	<i>CFG</i>	99.50	96.76	96.56	96.66
	<i>PDG</i>	98.55	83.55	90.01	86.66
	<i>Overall</i>	99.02	91.10	93.87	92.46

Table 1: Performance of NEURALPDA on **complete** code from Java and C/C++ (*Intrinsic Evaluation*)

<i>k</i>	F1-Score (in %)					
	<i>Java</i>			<i>C/C++</i>		
	<i>CFG</i>	<i>PDG</i>	<i>Overall</i>	<i>CFG</i>	<i>PDG</i>	<i>Overall</i>
3	99.70	93.24	97.17	98.39	91.10	96.01
4	99.41	92.51	96.61	98.33	90.18	95.28
5	99.26	91.23	95.96	97.91	89.10	94.37
6	99.04	90.36	95.40	97.14	87.91	93.31
7	98.75	89.45	94.82	96.69	86.45	92.39
8	98.44	88.70	94.29	96.66	86.66	92.46

Table 2: Performance of NEURALPDA on **partial** code from Java and C/C++ (*Intrinsic Evaluation*)

Empirical Evaluation (Extrinsic)

PDG from NEURALPDA



$$0 < VD\{PDG^*\} \leq VD\{PDG^\#\}$$



PDG from Program Analysis tool

Empirical Evaluation (Extrinsic)

PDG from NEURALPDA

$$0 < VD\{PDG^*\} \leq VD\{PDG^\#\}$$

PDG from Program Analysis tool

*The PDGs predicted by NEURALPDA approximates the performance of those generated by program analysis tools for **vulnerability detection on complete code** by 98.98%.*

Empirical Evaluation (Extrinsic)

PDG from NEURALPDA

$$0 < VD\{PDG^*\} \leq VD\{PDG^\#\}$$

PDG from Program Analysis tool

*The PDGs predicted by NEURALPDA approximates the performance of those generated by program analysis tools for **vulnerability detection on complete code** by 98.98%.*

*For the **vulnerability detection task on partial code**, the PDGs predicted by NEURALPDA helps an automated tool discover 14 real-world vulnerable code fragments.*

Empirical Evaluation (Qualitative)

Table 3: Performance of NEURALPDA on **different types of CFG/PDG edges** for Java (left) and C/C++ (right) code.

Graph	Edge Type	%C
CFG	<i>sequential</i>	99.54
	<i>if-else</i>	95.52
PDG	<i>data dependence</i>	82.78
	<i>control dependence</i>	96.33

Graph	Edge Type	%C
CFG	<i>sequential</i>	98.91
	<i>if-else</i>	**
PDG	<i>data dependence</i>	88.21
	<i>control dependence</i>	94.65

Building CFG/PDG with NEURALPDA: An Illustration

```
1 private boolean isValidUntil (Until annotation) {  
2     if (annotation != null) {  
3         double annotationVersion = annotation.value();  
4         if (annotationVersion <= version) {  
5             return false;  
6         }  
7     }  
8     return true;  
9 }
```

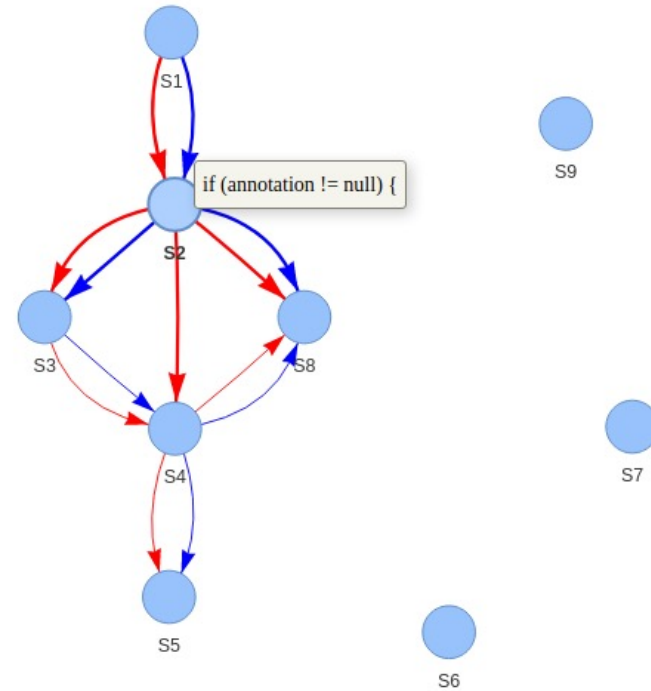


Figure 1: Java code listing (left) and its corresponding CFG/PDG (right) predicted by NEURALPDA

Key Takeaways

- ❖ NEURALPDA is the first neural network tool to predict program dependencies in complete as well as partial programs, which are accurate as well as **380x** faster to generate.
 - ❖ This work leads to a direction for improving program analysis (PA) for partial programs by combining pattern learning-based approaches with top-down PA techniques.
-

Conclusion



Why is this useful?

- ❖ Some existing approaches build CFG/PDGs for partial code via *manual intervention* (e.g., by wrapping around method signature) in a best-effort manner, often at the cost of several misses.

```
1 string subTag(string s, string a, string b) {
2   std::string lower_s;
3   std::transform(s.begin(), s.end(), lower_s.begin(), ::
4     tolower);
5   std::transform(a.begin(), a.end(), a.begin(), ::tolower);
6   auto position = lower_s.find(a);
   ...
}
```

In the above code listing, due to the unknown data type `string` on line 2 and API element `transform` on lines 3 and 4, traditional program analysis tools ignore all CFG/PDG edges to/from these statements. However, this listing is vulnerable on line 3, automated detection of which is not possible due to the missing edges.

Conclusion



Why is this useful?

- ❖ Some existing approaches build CFG/PDGs for partial code via *manual intervention* (e.g., by wrapping around method signature) in a best-effort manner, often at the cost of several misses.

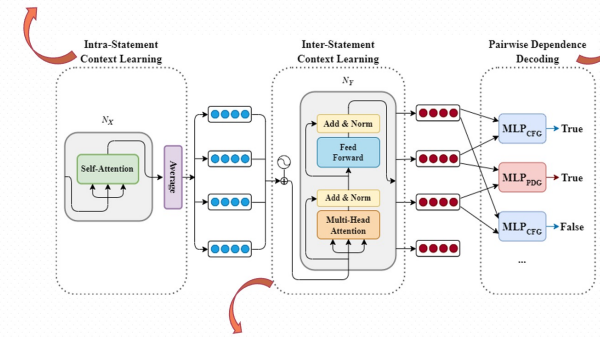
```
1 string subTag(string s, string a, string b) {  
2   std::string lower_s;  
3   std::transform(s.begin(), s.end(), lower_s.begin(), ::  
4     tolower);  
5   std::transform(a.begin(), a.end(), a.begin(), ::tolower);  
6   auto position = lower_s.find(a);  
   ...  
}
```

In the above code listing, due to the unknown data type `string` on line 2 and API element `transform` on lines 3 and 4, traditional program analysis tools ignore all CFG/PDG edges to/from these statements. However, this listing is vulnerable on line 3, automated detection of which is not possible due to the missing edges.

NEURALPDA: Neural Network-Based Program Dependence Analysis

Helps *relay the syntactic and semantic relationships* between the tokens in a statement to other statements in the code snippet.

The combination of all the CFG/PDG edges extracted via such an arc-factored approach is *realized as the CFG/PDG* for the given program.



The goal of this component is to *learn latent representations* for each statement that model the inter-statement dependencies.



Why is this useful?

- ❖ Some existing approaches build CFG/PDGs for partial code via *manual intervention* (e.g., by wrapping around method signature) in a best-effort manner, often at the cost of several misses.

```

1 string subTag(string s, string a, string b) {
2   std::string lower_s;
3   std::transform(s.begin(), s.end(), lower_s.begin(), ::
4     tolower);
5   std::transform(a.begin(), a.end(), a.begin(), ::tolower);
6   auto position = lower_s.find(a);
   ...

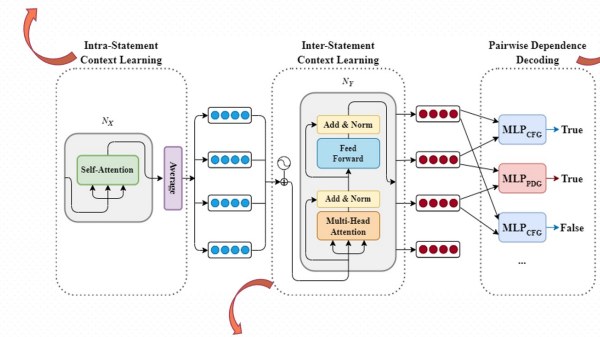
```

In the above code listing, due to the unknown data type `string` on line 2 and API element `transform` on lines 3 and 4, traditional program analysis tools ignore all CFG/PDG edges to/from these statements. However, this listing is vulnerable on line 3, automated detection of which is not possible due to the missing edges.

NEURALPDA: Neural Network-Based Program Dependence Analysis

Helps *relay the syntactic and semantic relationships* between the arc-tokens in a statement to other statements in the code snippet.

The combination of all the CFG/PDG edges extracted via such an arc-factored approach is *realized as the CFG/PDG* for the given program.



The goal of this component is to *learn latent representations* for each statement that model the inter-statement dependencies.

Empirical Evaluation

P/L	Graph	Accuracy	Precision	Recall	F1-Score
Java	CFG	99.79	98.31	98.58	98.44
	PDG	98.87	89.89	87.53	88.70
	Overall	99.33	94.75	93.83	94.29
C/C++	CFG	99.50	96.76	96.56	96.66
	PDG	98.55	83.55	90.01	86.66
	Overall	99.02	91.10	93.87	92.46

Table 1: Performance of NEURALPDA on *complete* code from Java and C/C++ (*Intrinsic Evaluation*)

k	F1-Score (in %)					
	Java			C/C++		
	CFG	PDG	Overall	CFG	PDG	Overall
3	99.70	93.24	97.17	98.39	91.10	96.01
4	99.41	92.51	96.61	98.33	90.18	95.28
5	99.26	91.23	95.96	97.91	89.10	94.37
6	99.04	90.36	95.40	97.14	87.91	93.31
7	98.75	89.45	94.82	96.69	86.45	92.39
8	98.44	88.70	94.29	96.66	86.66	92.46

Table 2: Performance of NEURALPDA on *partial* code from Java and C/C++ (*Intrinsic Evaluation*)

The PDGs predicted by NEURALPDA approximates the performance of those generated by program analysis tools for *vulnerability detection on complete code* by **98.98%** (*Extrinsic Evaluation*)

For the *vulnerability detection task on partial code*, the PDGs predicted by NEURALPDA helps an automated tool discover **14** real-world vulnerable code fragments (*Extrinsic Evaluation*)



Why is this useful?

- ❖ Some existing approaches build CFG/PDGs for partial code via *manual intervention* (e.g., by wrapping around method signature) in a best-effort manner, often at the cost of several misses.

```

1 string subTag(string s, string a, string b) {
2   std::string lower_s;
3   std::transform(s.begin(), s.end(), lower_s.begin(), ::
4     tolower);
5   std::transform(a.begin(), a.end(), a.begin(), ::tolower);
6   auto position = lower_s.find(a);
7   ...

```

In the above code listing, due to the unknown data type `string` on line 2 and API element `transform` on lines 3 and 4, traditional program analysis tools ignore all CFG/PDG edges to/from these statements. However, this listing is vulnerable on line 3, automated detection of which is not possible due to the missing edges.

Empirical Evaluation

P/L	Graph	Accuracy	Precision	Recall	F1-Score
Java	CFG	99.79	98.31	98.58	98.44
	PDG	98.87	89.89	87.53	88.70
	Overall	99.33	94.75	93.83	94.29
C/C++	CFG	99.50	96.76	96.56	96.66
	PDG	98.55	83.55	90.01	86.66
	Overall	99.02	91.10	93.87	92.46

Table 1: Performance of NEURALPDA on **complete** code from Java and C/C++ (*Intrinsic Evaluation*)

The PDGs predicted by NEURALPDA approximates the performance of those generated by program analysis tools for **vulnerability detection on complete code** by **98.98%** (*Extrinsic Evaluation*)

For the **vulnerability detection task on partial code**, the PDGs predicted by NEURALPDA helps an automated tool discover **14** real-world vulnerable code fragments (*Extrinsic Evaluation*)

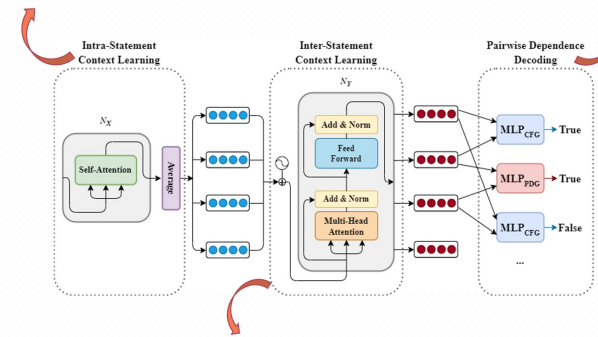
k	F1-Score (in %)					
	Java			C/C++		
	CFG	PDG	Overall	CFG	PDG	Overall
3	99.70	93.24	97.17	98.39	91.10	96.01
4	99.41	92.51	96.61	98.33	90.18	95.28
5	99.26	91.23	95.96	97.91	89.10	94.37
6	99.04	90.36	95.40	97.14	87.91	93.31
7	98.75	89.45	94.82	96.69	86.45	92.39
8	98.44	88.70	94.29	96.66	86.66	92.46

Table 2: Performance of NEURALPDA on **partial** code from Java and C/C++ (*Intrinsic Evaluation*)

NEURALPDA: Neural Network-Based Program Dependence Analysis

Helps *relay the syntactic and semantic relationships* between the arc-tokens in a statement to other statements in the code snippet.

The combination of all the CFG/PDG edges extracted via such an arc-factored approach is *realized as the CFG/PDG* for the given program.



The goal of this component is to *learn latent representations* for each statement that model the inter-statement dependencies.

Key Takeaways

- ❖ NEURALPDA is the first neural network tool to predict program dependencies in complete as well as partial programs, which are accurate as well as **380x** faster to generate.
- ❖ This work leads to a direction for improving program analysis (PA) for partial programs by combining pattern learning-based approaches with top-down PA techniques.



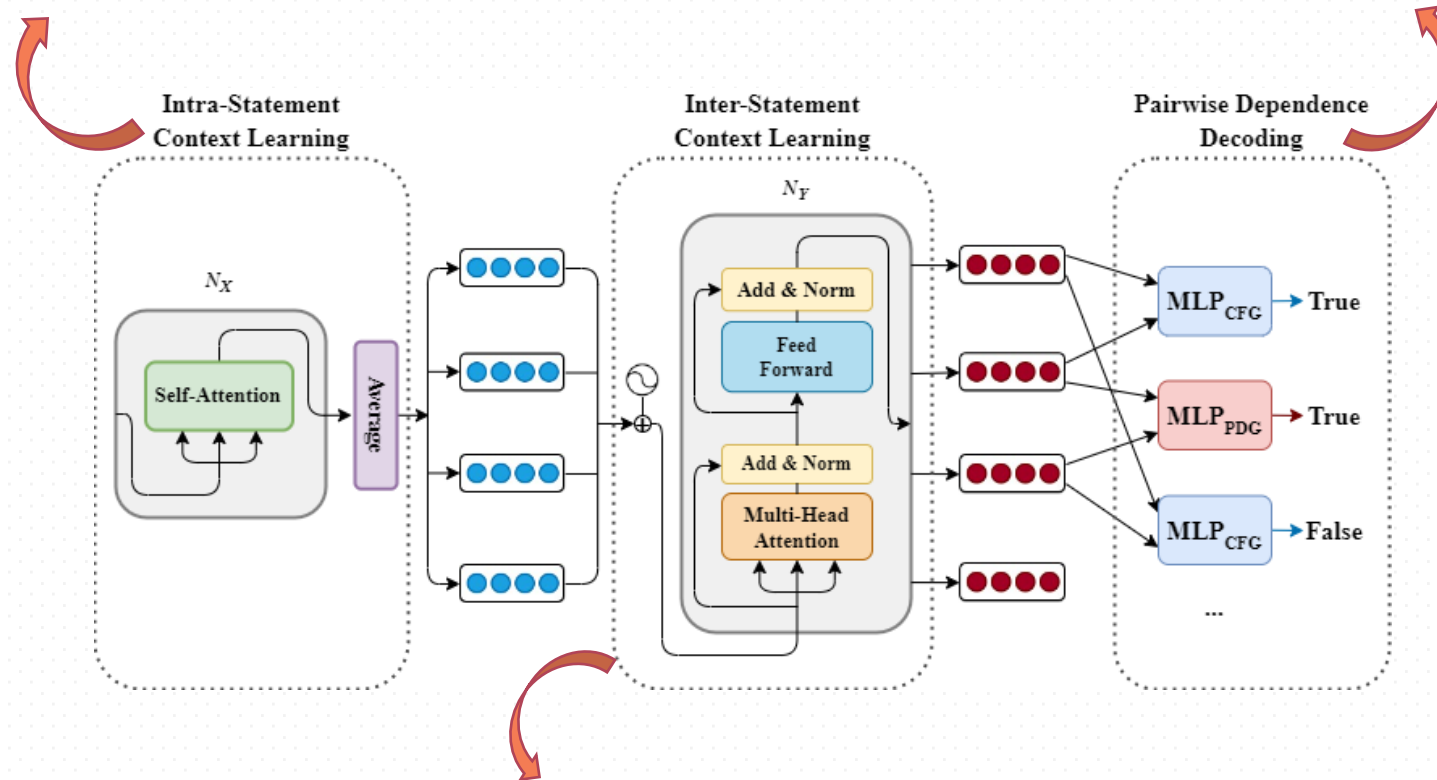
Conclusion

EXTRA SLIDES

NEURALPDA: Neural Network-Based Program Dependence Analysis

Helps *relay the syntactic and semantic relationships* between the tokens in a statement to other statements in the code snippet.

The combination of all the CFG/PDG edges extracted via such an arc-factored approach is *realized as the CFG/PDG* for the given program.



The goal of this component is to *learn latent representations* for each statement that model the inter-statement dependencies.

Building CFG/PDG with NEURALPDA: An Illustration

```
1 List<Game> allGames = gameMapper.getAllGamesByLeague(league);  
2 for (Game game : allGames) {  
3     game.getTeam1().setGame(game);  
4     game.getTeam2().setGame(game);  
5 }  
6 Collections.sort(allGames, new GameComparator());
```

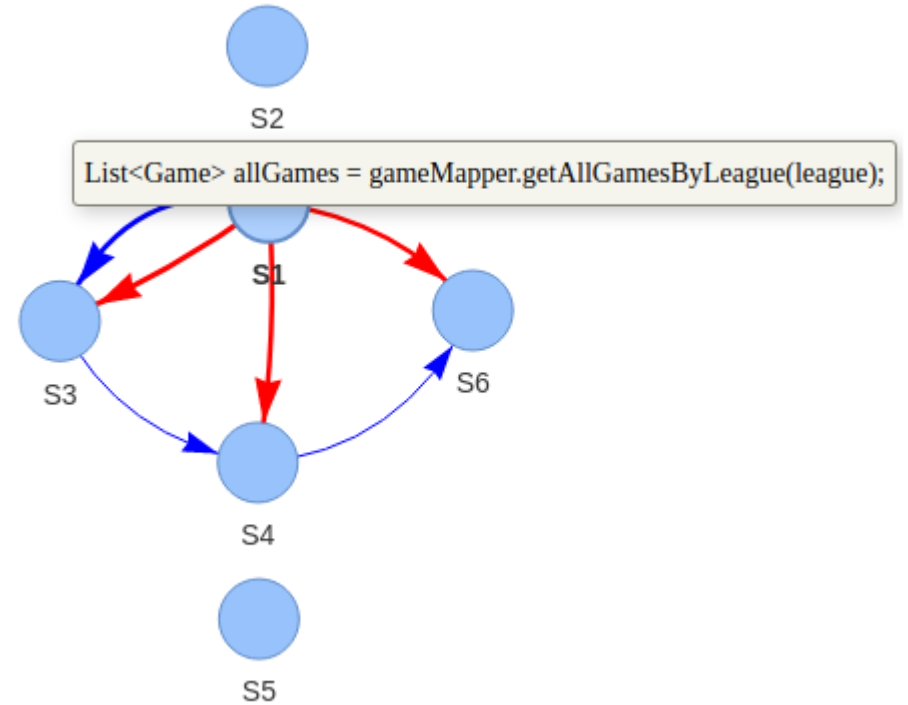


Figure 1: Partial Java code listing (left) and its corresponding CFG/PDG (right) predicted by NEURALPDA

Empirical Evaluation

P/L	Graph	Accuracy	Precision	Recall	F1-Score
Java	CFG	99.79	98.31	98.58	98.44
	PDG	98.87	89.89	87.53	88.70
	Overall	99.33	94.75	93.83	94.29
C/C++	CFG	99.50	96.76	96.56	96.66
	PDG	98.55	83.55	90.01	86.66
	Overall	99.02	91.10	93.87	92.46

Table 1: Performance of NEURALPDA on **complete** code from Java and C/C++ (*Intrinsic Evaluation*)

k	F1-Score (in %)					
	Java			C/C++		
	CFG	PDG	Overall	CFG	PDG	Overall
3	99.70	93.24	97.17	98.39	91.10	96.01
4	99.41	92.51	96.61	98.33	90.18	95.28
5	99.26	91.23	95.96	97.91	89.10	94.37
6	99.04	90.36	95.40	97.14	87.91	93.31
7	98.75	89.45	94.82	96.69	86.45	92.39
8	98.44	88.70	94.29	96.66	86.66	92.46

Table 2: Performance of NEURALPDA on **partial** code from Java and C/C++ (*Intrinsic Evaluation*)

The PDGs predicted by NEURALPDA approximates the performance of those generated by program analysis tools for **vulnerability detection on complete code** by **98.98%** (*Extrinsic Evaluation*)

For the **vulnerability detection task on partial code**, the PDGs predicted by NEURALPDA helps an automated tool discover **14** real-world vulnerable code fragments (*Extrinsic Evaluation*)

Conclusion

- ❖ NEURALPDA is the first neural network tool to predict program dependencies in complete as well as partial programs, which are accurate as well as **380×** faster to generate.
- ❖ This work leads to a direction for improving program analysis (PA) for partial programs by combining pattern learning-based approaches with top-down PA techniques.



All code, data, and supplementary material are available through this QR code.