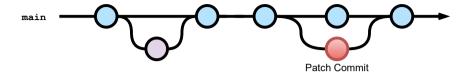
PPT4J: Patch Presence Test for Java Binaries

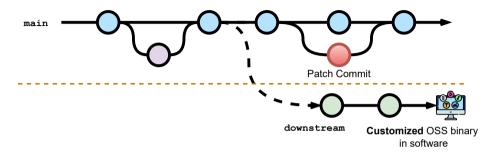
Zhiyuan Pan † Xing Hu *† Xin Xia † David Lo ‡ Xiaohu Yang †

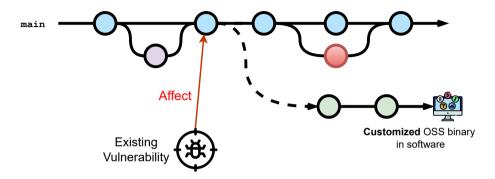
[†]The State Key Laboratory of Blockchain and Data Security, Zhejiang University

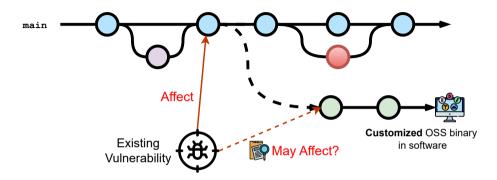
[‡]School of Computing and Information Systems, Singapore Management University

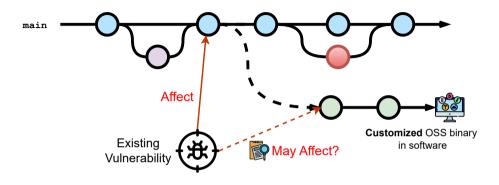
April 17, 2024











Patch Presence Test

Checks if a specific patch is applied to an unknown target binary.

Motivation

```
- YmlPullParser
                           parser = Xml.newPullParser():

    XPPAttributesWrapper attributes = new XPPAttributesWrapper(parser);

   + try
5
       XmlPullParser
                             parser = Xml.newPullParser():
       XPPAttributesWrapper attributes = new XPPAttributesWrapper(parser):
```

```
- Document<T> doc = parser.parse(is);
+ XMLStreamReader reader = StaxUtils.createXMLStreamReader(is):
+ Document<T> doc = parser_parse(reader):
```

```
if (A == Algorithm.none && B == 2 && C == 0) {
          return Mapper.deserialize(base64Decode(...). JWT.class):
        if (B == 2 && C == 0) {
          if (A == Algorithm.none) {
            return Mapper.deserialize(base64Decode(...). JWT.class):
6
            throw new InvalidJWTSignatureException()
```

Figure 1: Text diffs $^1 \neq$ Semantic changes



- Some (-) and (+) diff lines end up with no semantic changes.
- These lines introduce unrelated information to existing work that utilizes the complete patch diff.



CVE-2017-1000498, CVE-2016-8739, and CVE-2018-11797 respectively. Diffs are simplified for illustration.

Motivation

```
- YmlPullParser
                       parser = Xml.newPullParser():
- XPPAttributesWrapper attributes = new XPPAttributesWrapper(parser):
+ try
                         parser = Xml.newPullParser();
    XmlPullParser
    XPPAttributesWrapper attributes = new XPPAttributesWrapper(parser):
```

```
- Document<T> doc = parser.parse(is);
+ XMLStreamReader reader = StaxUtils.createXMLStreamReader(is):
+ Document<T> doc = parser_parse(reader):
```

```
if (A == Algorithm.none && B == 2 && C == 0) {
          return Mapper.deserialize(base64Decode(...). JWT.class):
        if (B == 2 && C == 0) {
          if (A == Algorithm.none) {
5
            return Mapper.deserialize(base64Decode(...). JWT.class):
6
            throw new InvalidJWTSignatureException()
```

Figure 1: Text diffs $^1 \neq$ Semantic changes



- Some (-) and (+) diff lines end up with no semantic changes.
- These lines introduce unrelated information to existing work that utilizes the complete patch diff.



 To extract precise semantic changes from diff that reflect all semantic information. while not including unrelated information.





Motivation

```
- YmlPullParser
                       parser = Xml.newPullParser():

    XPPAttributesWrapper attributes = new XPPAttributesWrapper(parser);

   XmlPullParser
                         parser = Xml.newPullParser():
   XPPAttributesWrapper attributes = new XPPAttributesWrapper(parser):
```

```
- Document<T> doc = parser.parse(is);
+ XMLStreamReader reader = StaxUtils.createXMLStreamReader(is):
+ Document<T> doc = parser_parse(reader):
```

```
if (A == Algorithm.none && B == 2 && C == 0) {
          return Mapper.deserialize(base64Decode(...). JWT.class):
        if (B == 2 && C == 0) {
          if (A == Algorithm.none) {
            return Mapper.deserialize(base64Decode(...). JWT.class):
6
            throw new InvalidJWTSignatureException()
```

Figure 1: Text diffs $^1 \neq$ Semantic changes



Facts

- Some (-) and (+) diff lines end up with no semantic changes.
- These lines introduce unrelated information to existing work that utilizes the complete patch diff.



 To extract precise semantic changes from diff that reflect all semantic information. while not including unrelated information.



Our Proposal

 A feature-based approach that highlights semantic changes.

CVE-2017-1000498. CVE-2016-8739. and CVE-2018-11797 respectively. Diffs are simplified for illustration.

Our Approach

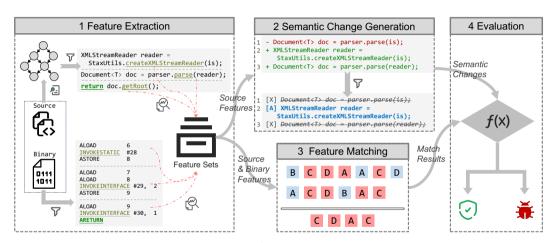


Figure 2: Overview of $\ensuremath{\mathrm{PPT4J}}$





Bytecode

Figure 3: Feature Extraction each generates a list of unified feature sets for source code and binary.

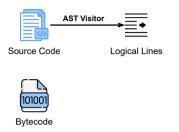


Figure 3: Feature Extraction each generates a list of unified feature sets for source code and binary.

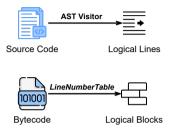


Figure 3: Feature Extraction each generates a list of unified feature sets for source code and binary.

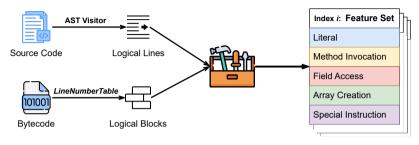


Figure 3: Feature Extraction each generates a list of unified feature sets for source code and binary.

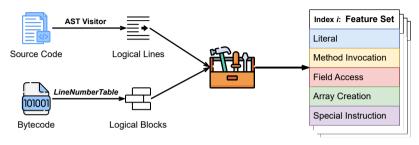


Figure 3: Feature Extraction each generates a list of unified feature sets for source code and binary.

- Literal ← Extract directly & Simplify arithmetic expressions
- Method ← Use type analysis to acquire signatures with more precise argument types

Special

- Distinctive operators: instanceof, ++/--, shift, ...
- Control flow manipulations: return, throw, if, loop¹ ···
- Syntactic sugars



¹We analyze CFGs of bytecode to distinguish condition & loop blocks.

Semantic Change Generation

A sliding window-based heuristic algorithm, which utilizes <u>feature set similarity</u> to filter out semantic redundant lines.

$$\mathcal{J}(A,B) = egin{cases} 1 & ext{A and B are both empty} \ rac{|A\cap B|}{|A\cup B|} & ext{otherwise} \end{cases}$$

- Split diff hunks into finer-grained blocks (Type-A, Type-D or Type-M).
- ② For each M-block, find the "optimal overlay" of (-) and (+) part.
- Sevaluate the similarity of feature sets within the overlay parts:
 - $0 \mathcal{J} = 1 \Rightarrow \text{mark as } excluded \text{ lines.}$
 - \mathfrak{O} > threshold $\sigma_f \Rightarrow$ mark as modification lines.
- Keep non-overlay parts as is, i.e., addition & deletion lines

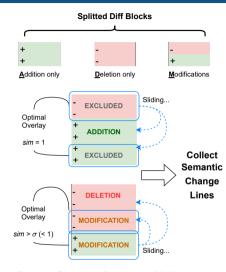


Figure 4: Examples of processing M-blocks

Feature Matching & Patch Presence Evaluation

What we have now:

- Lists of unified feature sets for:
 - Reference source code before patch
 - Reference source code after patch
 - m Target binary from user input
- Semantic change lines of the patch

It's time to figure out to what extent the binary resembles the diff part of reference sources.

• We apply the Longest Common Subsequence algorithm¹ to match the feature set sequences of the binary and the source code (i.e., (1), 11) and (1), 11).

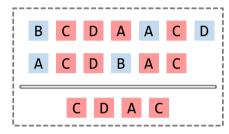


Figure 5: A simple example of the LCS algorithm. Each element in the sequence is a feature set.

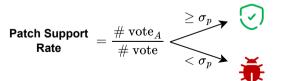
 $^{^1\}text{To}$ make the algorithm work, we define the equivalence of two elements (i.e., feature sets) as: $\mathcal{J}(A,B)>\sigma_f$

Feature Matching & Patch Presence Evaluation

What we have now:

- Lists of unified feature sets for:
 - Reference source code before patch
 - Reference source code after patch
 - m Target binary from user input
- Semantic change lines of the patch
- **NEW:** The matching results of (1), (11) and (11), (111)
- Each semantic change line votes for the final result: A) patched; B) unpatched Weighted vote: # votes = # features

- Addition line: if appears in Match(11), 11), vote **A**; otherwise vote **B**
- Deletion line: if not appears in Match(1),
 m), vote A; otherwise vote B
- Modification line pair (pre, post): if
 "J(post, binary) > J(pre, binary)", vote
 A; otherwise vote B
- Excluded lines are ignored in this procedure



Experimental Setup



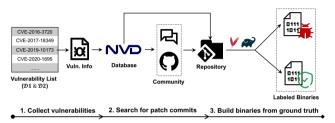


Figure 6: Steps to construct the dataset

- D1: The Java library vulnerabilities evaluated by the baseline
- D2: Vulnerabilities collected by Vul4J¹



 BScout² (reimplemented): A patch presence test framework specifically designed for Java binaries.



Metrics

- Accuracy
- Precision
- Recall
- F1 Score

¹Bui et al., "Vul4J: A Dataset of Reproducible Java Vulnerabilities Geared Towards the Study of Program Repair Techniques"

²Dai et al.. "BScout: Direct Whole Patch Presence Test for Java Executables"

Effectiveness

RQ.1: How accurate is the patch presence test framework compared to previous work?

Table 1: Test results on the dataset

Test Suite		Metrics			
		Accuracy	Precision	Recall	F1
BScout ¹	\mathcal{D} 1 2 \mathcal{D} 2	100% 87.9%	100% 100%	100% 75.8%	100% 86.2%
<u>Ррт4Ј</u>	$\mathcal{D}1$ $\mathcal{D}2$	100% 98.5%	100% 100%	100% 97.0%	100% 98.5%

- PPT4J does not generate false positive results.
- \bullet $\rm PPT4J$ outperforms the baseline BScout by 14.2% in terms of F1 score.



¹This refers to our reimplemented version.

 $^{^2\}text{The}$ results on $\boldsymbol{\mathcal{D}1}$ is consistent with the original paper.

RQ.1: How accurate is the patch presence test framework compared to previous work?

Table 1: Test results on the dataset

Test Suite		Metrics			
		Accuracy	Precision	Recall	F1
BScout ¹	\mathcal{D} 1 2 \mathcal{D} 2	100% 87.9%	100% 100%	100% 75.8%	100% 86.2%
<u>Ррт4Ј</u>	$\mathcal{D}1$ $\mathcal{D}2$	100% 98.5%	100% 100%	100% 97.0%	100% 98.5%

- PPT4J does not generate false positive results.
- PPT4J outperforms the baseline BScout by 14.2% in terms of F1 score.
- PPT4J is also effective in handling patches with minor changes.



¹This refers to our reimplemented version.

 $^{^2\}mathsf{The}$ results on $\boldsymbol{\mathcal{D}1}$ is consistent with the original paper.

Case Study: Minor Changes

Figure 7: CVE-2017-18349

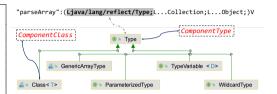


Figure 8: Class hierarchy of java.lang.reflect.Type

Case Study: Minor Changes

"parseArray":(Ljava/lang/reflect/Type; L...Collection; L...Object;)V

ComponentClass

GenericArrayType

Type

TypeVariable < D>

WildcardType

WildcardType

Figure 8: Class hierarchy of java.lang.reflect.Type

Figure 7: CVE-2017-18349

- parseArray(3) + parseArray(3) V.S.

BScout

- parseArray(Class, ...)

+ parseArray(Type, ...)

Ppt4J 😄

Efficiency

RQ.2: How **efficient** is the patch presence test framework, especially when dealing with large code repositories?

Table 2: Time consumption¹ on the dataset

Framework	Average	~75%
BScout	0.34 sec/patch	0.28 sec/patch
Ppt4J	0.48 sec/patch	0.30 sec/patch

- Most patches can be quickly analyzed.
- Time cost is not proportional to the project size because only dependent² bytecodes are analyzed.
- A bit slower than BScout, but the advantages in effectiveness can compensate for this.



¹The startup time of the virtual machine and third-party dependencies is not considered.

²Java classes fixed by the patch, and their dependent classes.

Ablation Study

RQ.3: How do the analyses in Feature Extraction¹ contribute to the overall effectiveness?

Four variants:

- PPT4J_**FULL**: Complete version
- **2** PPT $4J_\Delta 1$: Remove type analysis
- **3** PPT4J $_{\Delta}$ 2: Ignore special instructions (e.g., loop and branch)
- **4** PPT4J $_{\Delta}$ 3: Remove constant propagation/folding



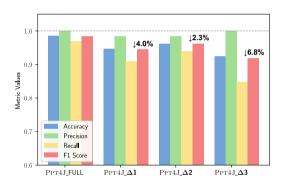


Figure 9: Test results for different variants of PPT4J



¹Please refer to Section 3.2.3 of our paper for more details.

Ablation Study

RQ.3: How do the analyses in Feature Extraction¹ contribute to the overall effectiveness?

Four variants:

- PPT4J_**FULL**: Complete version
- **2** PPT $4J_{\Delta}1$: Remove type analysis
- PPT4J_Δ2: Ignore special instructions (e.g., loop and branch)
- **4** PPT4J $_{\Delta}$ 3: Remove constant propagation/folding

The above analyses contribute to the performance improvement.

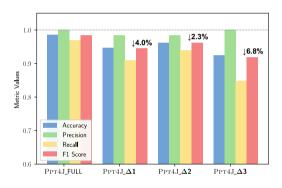


Figure 9: Test results for different variants of PPT4J



¹Please refer to Section 3.2.3 of our paper for more details.

In-the-wild Evaluation

RQ.4: Can our approach analyze open-source libraries in real-world applications?

Test Procedures:

- Extract OSS binaries from IntelliJ IDEA¹.
- 2 Run PPT² with specific patches.
- 3 Utilize unit tests in patch commits.
- 4 Check if UT results = PPT results.

Table 3: Results on open-source libraries within IntelliJ IDEA

.=.	• •	Version Timeline				
涆	وا	V1	V2	VЗ	V4	V 5
CVE-2019-12402	08/19	TN	TN	TP	TP	TP
CVE-Anonymous-1	-	TN	TN	TN	TN	TN
CVE-Anonymous-2	-	TN	TN	TN	TN	TN
CVE-2021-29425	05/18	TN	TN	TP	TP	TP
HTTPCLIENT-1803	01/17	TN	FN	FN	FN	FN
CVE-2017-1000487	10/13	TP	TP	TP	TP	TP
CVE-2015-6748	07/15	N/A	TP	TP	TP	TP
CVE-2015-6420	11/15	TN	TP	TP	TP	TP



¹https://www.jetbrains.com/idea

²Patch Presence Test

In-the-wild Evaluation

RQ.4: Can our approach analyze open-source libraries in real-world applications?

Test Procedures:

- Extract OSS binaries from IntelliJ IDEA¹.
- 2 Run PPT² with specific patches.
- 3 Utilize unit tests in patch commits.
- Oheck if UT results = PPT results.
- PPT4J achieves 89.7% accuracy with no false positive results.
 (BScout accuracy: 76.9%, ↓ 14.3%)
- PPT4J detects two un-patched vulnerabilities. We have reported this potential problem to the vendor.

Table 3: Results on open-source libraries within IntelliJ IDEA

بھ	••	Version Timel							
II	وا	V1	V2	VЗ	V4	Р ТР			
CVE-2019-12402	08/19	TN	TN	TP	TP	TP			
CVE-Anonymous-1	-	TN	TN	TN	TN	TN			
CVE-Anonymous-2	-	TN	TN	TN	TN	TN			
CVE-2021-29425	05/18	TN	TN	TP	TP	TP			
HTTPCLIENT-1803	01/17	TN	FN	FN	FN	FN			
CVE-2017-1000487	10/13	TP	TP	TP	TP	TP			
CVE-2015-6748	07/15	N/A	TP	TP	TP	TP			
CVE-2015-6420	11/15	TN	TP	TP	TP	TP			



¹https://www.jetbrains.com/idea

²Patch Presence Test

Conclusion

In summary, we made the following contributions:

- **1** A patch presence test framework for Java binaries: PPT4J.
- 2 A dataset of binaries to evaluate the effectiveness of PPT4J.
- $\ensuremath{\mathfrak{g}}$ We evaluate $\mathrm{PPT}4J,$ with results suggesting that $\mathrm{PPT}4J$ outperforms the baseline and is also capable in real-world scenarios.
- We release the replication package to facilitate future research.

Thanks for your attention!







