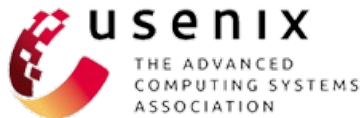# EaTVul: ChatGPT-based Evasion Attack Against Software Vulnerability Detection

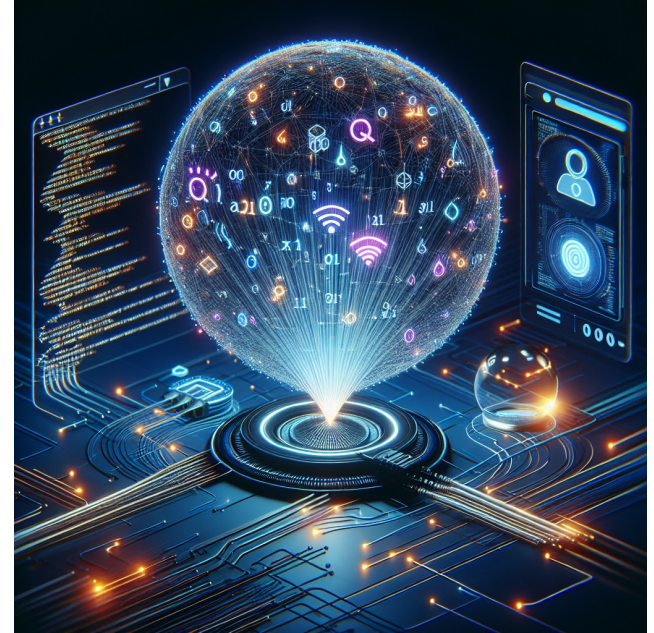**Shigang Liu**[1,2], Di Cao[2], Junae Kim[3], Tamas Abraham[3], Paul Montague[3], Seyit Camtepe[1], Jun Zhang[2], and Yang Xiang[2]

[1]CSIRO's Data61, [2]Swinburne University of Technology, [3]DST Group, Australia

# Background

- Machine learning-as-a-service has been widely applied in software security.

- Adversarial learning has long been a threat for cybersecurity.

- There is lack of thorough assessment the security issues when facing adversarial learning.

# A Motivation Example

# Our work aims to…

- Explore the susceptibility of machine learning/deep neural network models to adversarial attacks.

- Develop an effective scheme to generate adversarial code and inject it into vulnerable samples to bypass deep neural network systems.

# Assumption

- **Attacker's capability:** perturb the test queries and query the deployed vulnerability detection model.

- **Attacker's Knowledge:** no access to the architecture and parameters of target models.

- **Attacker's Goal:** deceive the targeted vulnerability detection tools through imperceptible modifications to the inputs.

# Proposed EaTVul



1 : step 1 - adversarial data generation    2 : step 2 — adversarial learning

# Proposed EaTVul



① : step 1 - adversarial data generation    ② : step 2 — adversarial learning

# Proposed EaTVul



① : step 1 - adversarial data generation    ② : step 2 – adversarial learning

# Proposed EaTVul

# Proposed EaTVul

# Example of Identified Features

- Display of the important features identified by attention mechanism, with the importance decreases from red to yellow.

- The identified features include *static, const, strstr, strchr, val, sscanf*.

```
1:  static void sdp_fmtp_get(const char *attributes, const char *name, int *attr)
2:  {
3:      const char *kvp = "";
4:      int val;
5:
6:      if (attributes && !(kvp = strstr(attributes, name))) {
7:      return;
8:      }
9:
10:    if (kvp != attributes && *(kvp - 1) != ' ' && *(kvp - 1) != ';') {
11:        /* Keep searching as it might still be in the attributes string */
12:        sdp_fmtp_get(strchr(kvp, ';'), name, attr);
13:    } else if (sscanf(kvp, "%*[^=]=%30d", &val) == 1) {
14:        *attr = val;
15:    }
16:  }
```

# Example of Generated Adversarial Data

```c
1   #include <stdio.h>
2
3   struct a_chan {
4       int value;
5   };
6
7   struct b_const {
8       double value;
9   };
10
11  static const char a_array[] = "NULL";
12
13  void myFunction() {
14
15      struct a_chan myAChan;
16      struct b_const myBConst;
17
18      myAChan.value = 0;
19      myBConst.value = 0.0;
20
21      for (int i = 0; i < 10; i++) {
22          printf("Iteration %d:\n", i);
23
24          myAChan.value += i;
25          myBConst.value += 0.5 * i;
26
27          printf("a_chan value: %d\n", myAChan.value)
28          printf("b_const value: %f\n", myBConst.value)
29      }
30
31  }
```

(a) Raw data generated by ChatGPT

```c
1   struct a_chan { int value; } queue;
2   struct b_const { double value; } p_project;
3   static const char a_array[] = "NULL";
4
5   queue.value = 0;
6   p_project.value = 0.0;
7
8   for (int i = 0; i < 10; i++)
9       printf("i = %d\n", i);
```

(b) Further optimized data by ChatGPT

CSIRO

# Visualization

- The distribution of vulnerable, non-vulnerable samples, and the adversarial samples.



(a) Vulnerable and non-vulnerable (normal) samples          (b) Movement of adversarial samples

# Evaluation

- RQ1: How effective is fuzzy genetic algorithm in selecting the seed adversarial data compared with randomization?

- RQ2: How effective is EaTVul based on adversarial data generated by ChatGPT originally and after optimization?

- RQ3: How effective is EaTVul with recently developed machine learning-based software vulnerability detection systems?

- RQ4: How effective is EaTVul when compared with state-of-the-art large language models (LLM) and other machine learning tools that using BiLSTM for software vulnerability detection?

- RQ5: How EaTVul behaves/performs regarding obfuscation/diversification methods?

- RQ6: How effective is EaTVul when generalized to other programming languages?

# RQ4: EaTVul attack BiLSTM & LLMs

- Conducted experiment based on the datasets of CWE119, CWE399, Asterisk and OpenSSL.

- The target models include Poster-Lin, MDVD, CodeBERT, and CodeGen.

| Target Model | Data | Snippet Size = 2 | | | | Snippet Size =3 | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Top@5 | Top@ 10 | Top@15 | Top@ 20 | Top@ 5 | Top@ 10 | Top@15 | Top@20 |
| Poster-Lin | Asterisk | 0.900 | 0.867 | 0.756 | 0.758 | 1.000 | 1.000 | 0.878 | 0.858 |
| | OpenSSL | 1.000 | 1.000 | 0.745 | 0.717 | 1.000 | 1.000 | 0.911 | 0.875 |
| | CWE119 | 0.933 | 0.934 | 0.899 | 0.867 | 1.000 | 0.967 | 0.956 | 0.925 |
| | CWE399 | 1.000 | 0.833 | 0.823 | 0.767 | 1.000 | 1.000 | 1.000 | 0.917 |
| MDVD | Asterisk | 1.000 | 0.867 | 0.844 | 0.784 | 1.000 | 1.000 | 1.000 | 0.975 |
| | OpenSSL | 1.000 | 1.000 | 0.845 | 0.817 | 1.000 | 1.000 | 1.000 | 1.000 |
| | CWE119 | 1.000 | 0.933 | 0.877 | 0.825 | 1.000 | 1.000 | 0.978 | 0.958 |
| | CWE399 | 1.000 | 0.899 | 0.867 | 0.767 | 1.000 | 1.000 | 0.967 | 0.950 |
| CodeBERT | Asterisk | 0.900 | 0.850 | 0.803 | 0.740 | 0.900 | 0.885 | 0.845 | 0.832 |
| | OpenSSL | 0.800 | 0.800 | 0.768 | 0.735 | 0.900 | 0.864 | 0.858 | 0.835 |
| | CWE119 | 0.900 | 0.840 | 0.834 | 0.785 | 1.000 | 0.935 | 0.911 | 0.865 |
| | CWE399 | 0.900 | 0.825 | 0.786 | 0.776 | 1.000 | 0.920 | 0.878 | 0.858 |
| CodeGen | Asterisk | 0.900 | 0.867 | 0.844 | 0.784 | 0.933 | 0.925 | 0.899 | 0.867 |
| | OpenSSL | 0.900 | 1.000 | 0.845 | 0.817 | 0.956 | 0.911 | 0.875 | 0.845 |
| | CWE119 | 1.000 | 0.933 | 0.877 | 0.825 | 1.000 | 1.000 | 0.928 | 0.875 |
| | CWE399 | 1.000 | 0.899 | 0.867 | 0.767 | 1.000 | 0.950 | 0.917 | 0.880 |

# RQ5: EaTVul attack Obfuscation/Diversification

- The baseline/attack models are EaTVul, Differentiable Obfuscator and Milo.

- The target models include Asteria and LineVul.

- Uses Juliet C/C++ Test Suite Datasets.



| Dataset | Target Model | Attack Model | ASR | F1-Score |
|---------|--------------|--------------|-----|----------|
| CWE119 | Asteria | Differentiable Obfuscator | 66.40 | 81.45 |
| | | Milo | 35.80 | |
| | | EaTVul | **99.50** | |
| | LineVul | Differentiable Obfuscator | 63.50 | 83.45 |
| | | Milo | 44.30 | |
| | | EaTVul | **92.30** | |
| CWE399 | Asteria | Differentiable Obfuscator | 58.80 | 82.60 |
| | | Milo | 26.30 | |
| | | EaTVul | **89.50** | |
| | LineVul | Differentiable Obfuscator | 62.80 | 83.50 |
| | | Milo | 28.70 | |
| | | EaTVul | **89.20** | |
| CWE416 | Asteria | Differentiable Obfuscator | 52.50 | 80.40 |
| | | Milo | 20.35 | |
| | | EaTVul | **84.30** | |
| | LineVul | Differentiable Obfuscator | 58.60 | 81.70 |
| | | Milo | 19.80 | |
| | | EaTVul | **87.50** | |

# RQ6: Generalization of EaTVul

- Evaluate on Java programs.

- The attack models are EaTVul, Differentiable Obfuscator and Milo.

- The target models include FUNDED and VDet.

| Target Model | Attack Model | ASR | F1-Score |
|---|---|---|---|
| FUNDED | Differentiable Obfuscator | 53.50 | 85.35 |
| | Milo | 42.70 | |
| | EaTVul | **88.60** | |
| VDet | Differentiable Obfuscator | 63.50 | 86.60 |
| | Milo | 62.80 | |
| | EaTVul | **87.30** | |

# Conclusion

- We propose a novel evasion attack approach, named EatVul, which produces code to evade ML/deep neural network-based vulnerability detectors.

- We have conducted an evaluation of EatVul against state- of-the-art baselines, and the experimental results demonstrated that our scheme achieved a 100% success rate in most cases with a snippet size of 4.

- We have made our proposed system, EatVul, available to the research community. The datasets and code are available at **https://github.com/wolong3385/EatVul-Resources**.

- This work demonstrates the susceptibility of vulnerable samples to manipulation and highlights the need for robust defense mechanisms capable of mitigating such adversarial attacks.

**Thank you!**
**Q&A**