

# CE 815 – Secure Software Systems

---

ML-Based Vulnerability Detection Methods (Hoppity)

Mohammad Haddadian/Mehdi Kharrazi  
Department of Computer Engineering  
Sharif University of Technology



Acknowledgments: Some of the slides are fully or partially obtained from other sources. A reference is noted on the bottom of each slide, when the content is fully obtained from another source. Otherwise a full list of references is provided on the last slide.



# Introduction

---

- Vulnerability detection as first step
- Then, Vulnerability repair

**HOPPITY: Learning Graph Transformations to Detect and Fix Bugs in Programs**, E. Dinella, H. Dai, Z. Li, M. Naik, L. Song, K. Wang, ICLR 2020.



# Problem

---

Source-code analysis is:

- Undecidable
- Noisy
- Rules are hand written
- Tailored to specific code bases / bug patterns



# Javascript Challenges

---

- Incorrect operators
- Incorrect identifiers
- Accessing undefined properties
- Mishandling variable scopes
- Type incompatibilities

# Example



```
function clearEmployeeListOnLinkClick() {
  document.querySelector("a").addEventListener("click",
    function(event) {
      document.querySelector("ul").InnerHTML = "";
    }
  );
}
```

(a) InnerHTML should have been innerHTML.

```
if (matches) {
  return {
    episode: Number(matches.groups.episode),
    hosts: matches.groups.hosts.split(/([, &]+|\sand\s)/).
      map(el => S(el).trim().s)
  };
}
```

(b) Highlighted parentheses should have been removed.

```
module.exports = function (grunt) {
  grunt.initConfig({
    execute: {...}, copy: {...}, checktextdomain: {...}
    wp_readme_to_markdown: {...}, makepot: {...}})
  ...
  grunt.registerTask('default', ['wp_readme_to_markdown',
    'makepot', 'execute', 'checktextdomain'])
};
```

(c) copy function should have also been included in the highlighted list.

```
export default {
  computed: {
    level () {
      return dictMap.skillLevel[
        parseInt((this.value === 0 ? 1 : this.value)/20)];
    }, ...
  }
}
```

(d) parseInt should have been removed because === implies this.value is an integer.

# Solution



---

Leverage large amounts of Javascript fixes on Github to locate and repair bugs

# Steps

---



- Represent source code
- Represent fixes
- Learning



# Goal

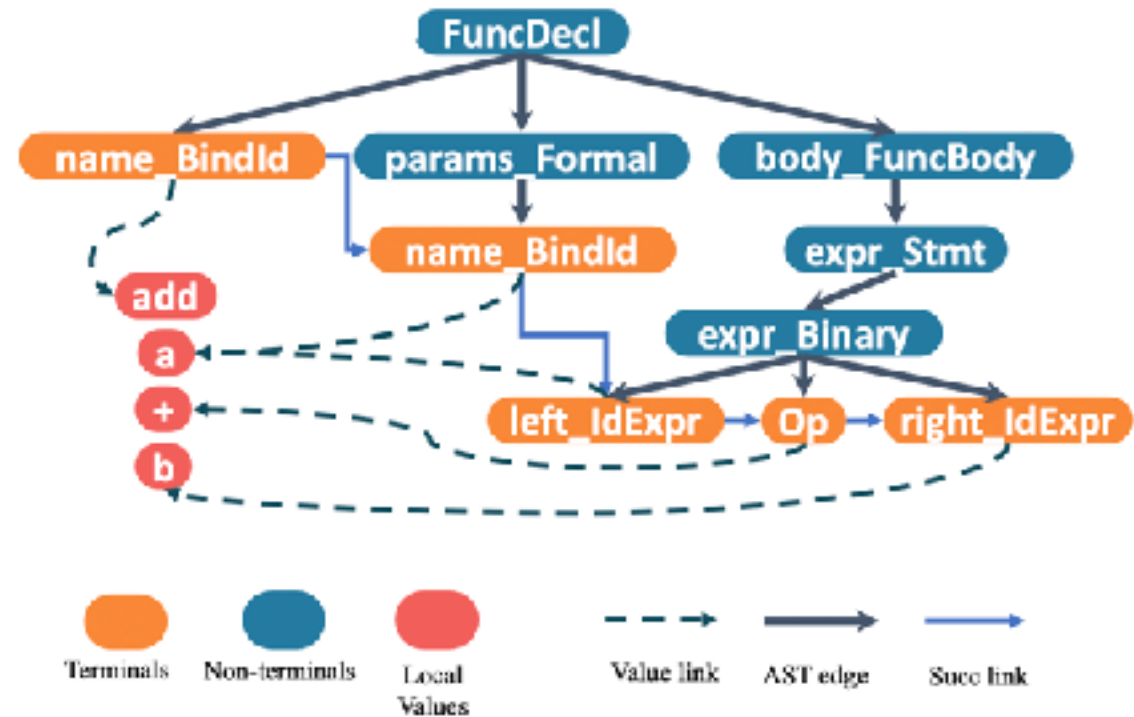


- 1- `function add(a) { a + b; }` Buggy
- 2- `function add(a, b) { a + b; }` Step1
- 3- `function add(a, b) { return a + b; }` Step2



# Source code representation

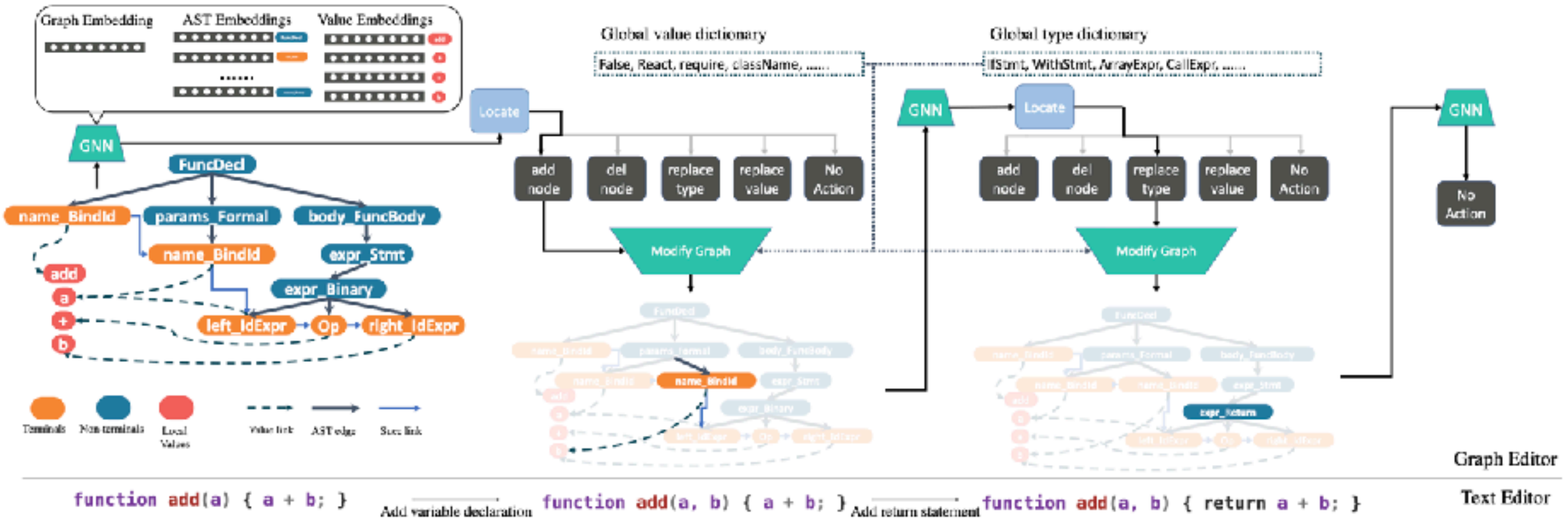
- AST
- ValueLink

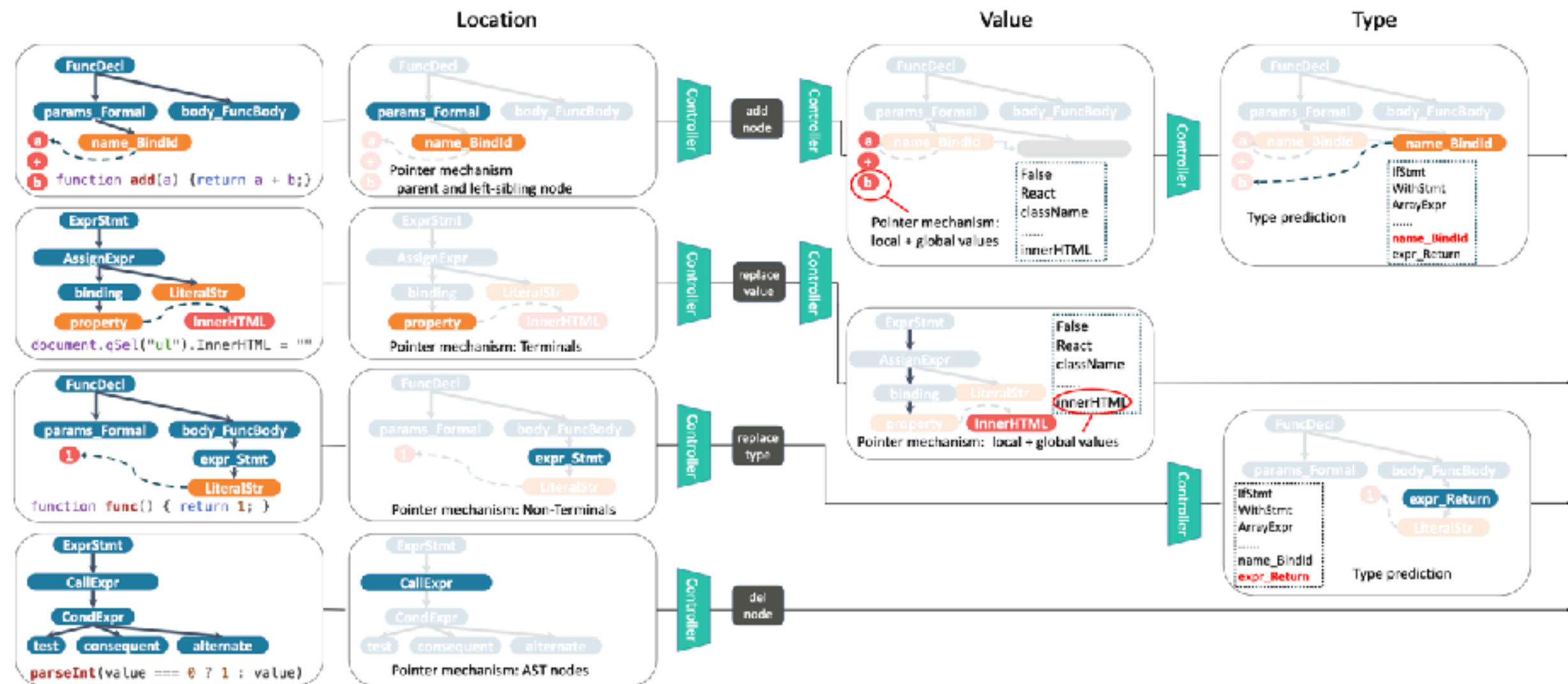




# Fix representation

- Graph Edits

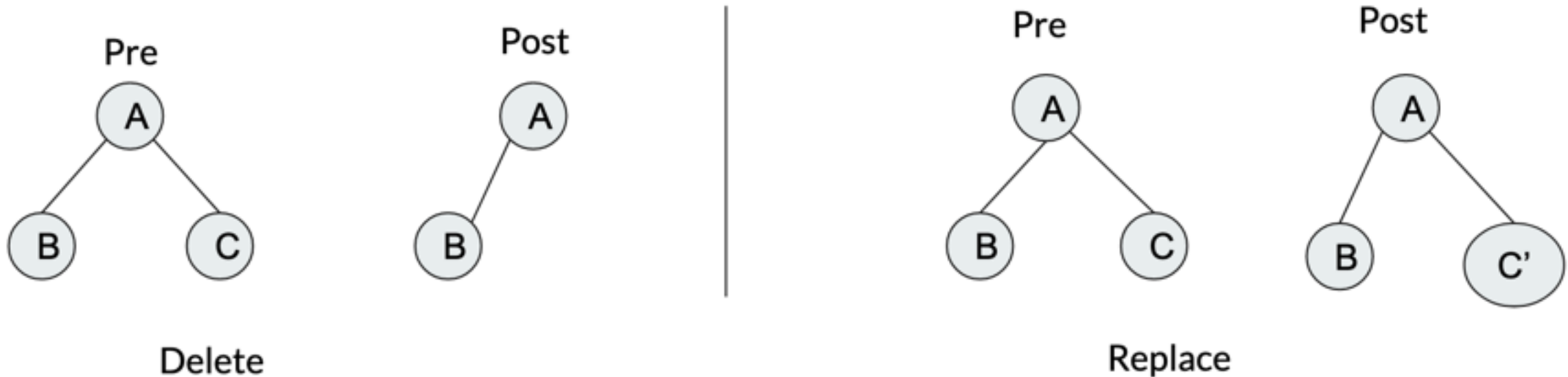






# Low level primitives

- Location
- Value
- Type



# Graph edit operators

---



# Graph transformation

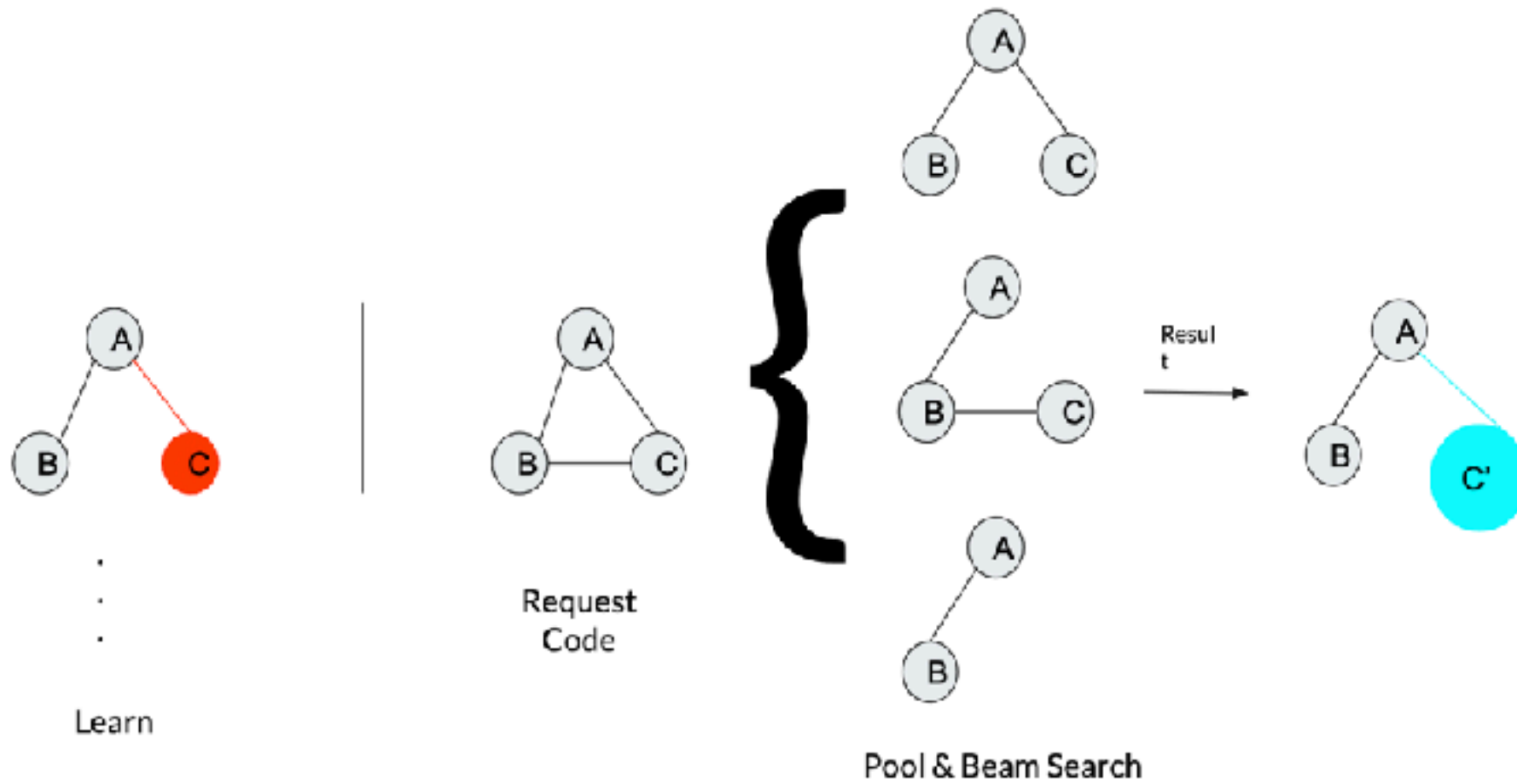


1

2

3

# Inference







# Dataset

- OneDiff (just one change)
- ZeroOneDiff (zero or one edit)
- ZeroOneTwoDiff (zero, one or two edits)

	ADD	REP_TYPE	REP_VAL	DEL	total
train	6,473	1,864	251,097	31,281	290,715
validate	790	245	31,357	3,957	36,349
test	796	233	31,387	3,945	36,361

**Table 1: Statistic of OneDiff dataset. See appendix for more information of other dataset.**

# Evaluation



	Total		<i>Location</i>		<i>Operator</i>	<i>Value</i>		<i>Type</i>	
	Top-3	Top-1	Top-3	Top-1	Top-1	Top-3	Top-1	Top-3	Top-1
<b>TOTAL</b>	<b>26.1</b>	14.2	35.5	20.4	34.4	52.3	29.1	76.1	66.7
ADD	52.9	39.2	69.6	51.4	70.6	65.7	55.1	76.8	68.5
REP_VAL	23.4	11.9	33.3	18.5	31.7	53.0	28.8	-	-
REP_TYPE	71.7	52.4	73.0	52.8	79.4	-	-	74.7	61.0
DEL	39.6	24.8	44.0	27.5	45.8	-	-	-	-
Random	.08	.07	2.28	1.4	27.7	.01	.01	.27	0

Table 2: Evaluation of model on the OneDiff dataset: accuracy (%).



# Evaluation (cont.)

Type	GGNN-Rep	GGNN-Cls	HOPPITY
Top-1	53.2%	<b>99.6%</b>	90.0%
Top-3	85.8%	<b>99.6%</b>	94.8%

Table 3: REP\_TYPE accuracies with location+op.

Value	GGNN-Rep	GGNN-RNN	HOPPITY
Top-1	63.8%	60.3%	<b>69.1%</b>
Top-3	67.6%	63.6%	<b>73.4%</b>

Table 4: REP\_VAL accuracies with location+op.

	Top-1	Top-3
HOPPITY	<b>67.7%</b>	<b>73.3%</b>
SequenceR	64.2%	68.6%

Table 5: Overall OneDiff accuracy with location.

Bug Type	Amount	TAJS	HOPPITY
Undefined Property	7	0	1
Functional Bug	11	0	3
Refactoring	12	0	1
Total	30	0	5

Table 6: Comparison with TAJS.

**Dos and Don'ts of Machine Learning in Computer Security**, D. Arp, E. Quiring, F. Pendlebury, A. Warnecke, F. Pierazzi, C. Wressnegger, L. Cavallaro, K. Rieck, Usenix 2022.

**Machine Learning already solved  
many problems in computer security**

Unfortunately not... 🙄



# Motivation—Historical Examples



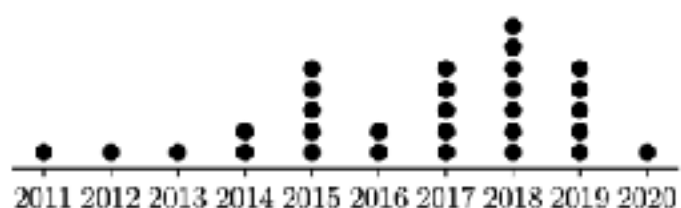
## **Network intrusion detection: The base rate fallacy**

- Intrusion detectors should have low false positive rates (FPR)
- 'Low' FPR often still corresponds to large number of false positives

## **Android malware detection: Spatio-temporal bias inflating performance**

- Models trained with access to 'future' information
- Unrealistic class balance inflates performance

# Overview



## 1. Identification of common pitfalls

- 10 subtle issues affecting ML for security
- Recommendations for avoiding them

## 2. Survey on the prevalence of pitfalls

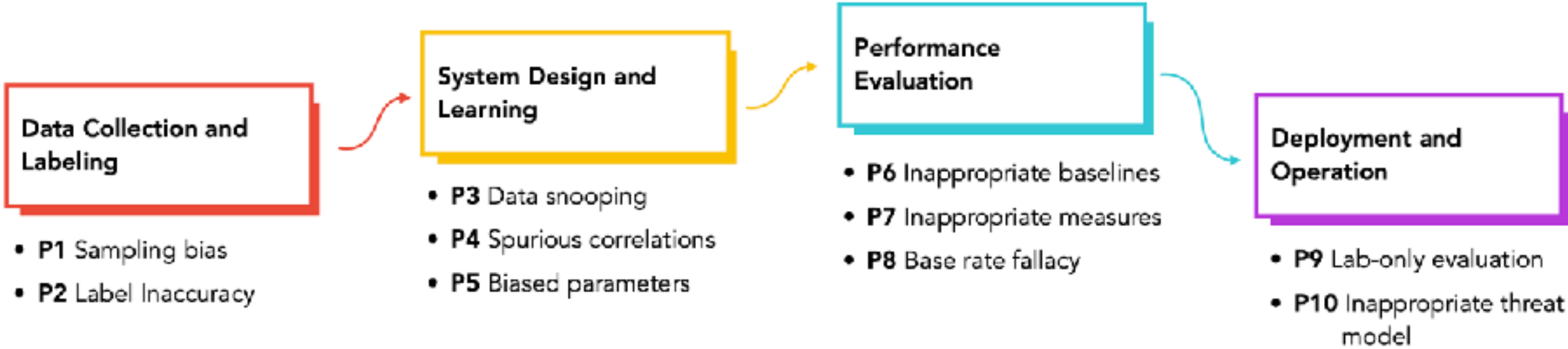
- Review of 30 top papers in security
- Pitfalls are widespread

## 3. Case studies demonstrating impact of pitfalls

- Mobile malware detection
- Vulnerability discovery
- Source code authorship attribution
- Network intrusion detection



# ML Pipeline and Pitfalls



**P1 – Sampling Bias.** The collected data does not sufficiently represent the true data distribution of the underlying security problem [1, 30, 33].

90% present

**P2 – Label Inaccuracy.** The ground-truth labels required for classification tasks are inaccurate, unstable, or erroneous, affecting the overall performance of a learning-based system [85, 144].

100% present

**P3 – Data Snooping.** A learning model is trained with data that is typically not available in practice. Data snooping can occur in many ways, some of which are very subtle and hard to identify [1].

50% present

**P4 – Spurious Correlations.** Artifacts unrelated to the security problem create shortcut patterns for separating classes. Consequently, the learning model adapts to these artifacts instead of solving the actual task.

100% present

**P5 – Biased Parameter Selection.** The final parameters of a learning-based method are not entirely fixed at training time. Instead, they indirectly depend on the test set.

100% present

**P6 – Inappropriate Baseline.** The evaluation is conducted without, or with limited, baseline methods. As a result, it is impossible to demonstrate improvements against the state of the art and other security mechanisms.

100% present

**P7 – Inappropriate Performance Measures.** The chosen performance measures do not account for the constraints of the application scenario, such as imbalanced data or the need to keep a low false-positive rate.

33% present

**P8 – Base Rate Fallacy.** A large class imbalance is ignored when interpreting the performance measures leading to an overestimation of performance.

100% present

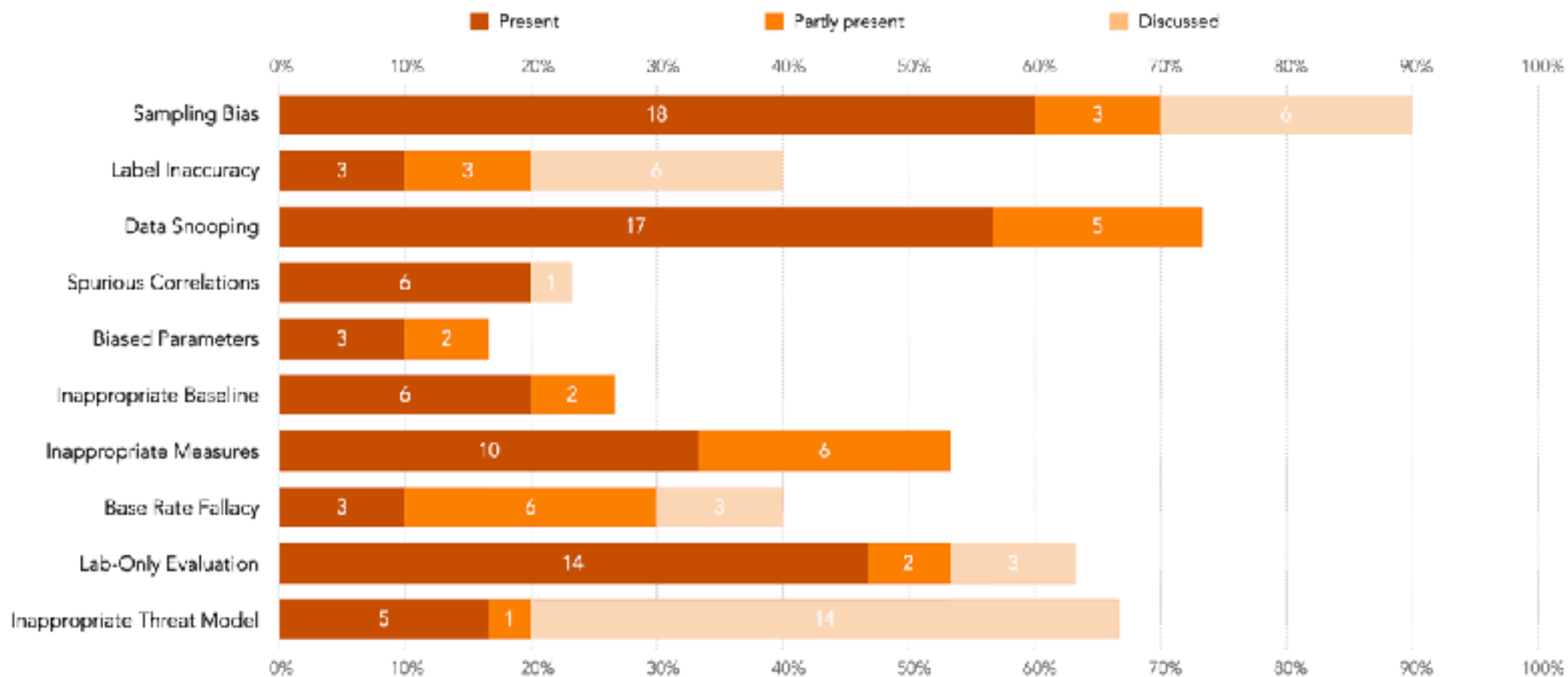
**P9 – Lab-Only Evaluation.** A learning-based system is solely evaluated in a laboratory setting, without discussing its practical limitations.

17% present

**P10 – Inappropriate Threat Model.** The security of machine learning is not considered, exposing the system to a variety of attacks, such as poisoning and evasion attacks.

17% present

# Prevalence Study



# Impact Analysis

## Android Malware Detection

- P1:** Sampling Bias
- P4:** Spurious Correlations
- P7:** Inappropriate Performance Measures

## Authorship Attribution

- P1:** Sampling Bias
- P4:** Spurious Correlations

## Vulnerability Discovery

- P2:** Label Inaccuracy
- P4:** Spurious Correlations
- P6:** Inappropriate Baselines

## Network Intrusion Detection

- P6:** Inappropriate baselines
- P9:** Lab-only evaluation

# Acknowledgments

---



- [Hoppity] HOPPITY: Learning Graph Transformations to Detect and Fix Bugs in Programs, E. Dinella, H. Dai, Z. Li, M. Naik, L. Song, K. Wang, ICLR 2020.
- [Arp] Dos and Don'ts of Machine Learning in Computer Security, D. Arp, E. Quiring, F. Pendlebury, A. Warnecke, F. Pierazzi, C. Wressnegger, L. Cavallaro, K. Rieck, Usenix 2022.