# CENTRIS: A Precise and Scalable Approach for Identifying Modified Open-Source Software Reuse

43rd International Conference on Software Engineering

**Seunghoon Woo\*,** Sunghan Park\*, Seulbae Kim$^{\dagger}$, Heejo Lee\*, Hakjoo Oh\*

**\*Korea University**, $^{\dagger}$Georgia Institute of Technology

ICSE 2021

# GOAL

- **Identifying Open-source software (OSS) components in the target software**

- **Motivation**
  - Open-source software is reused extensively in software development
  - Reusing OSS without <u>proper management</u>
    - 🙁 Vulnerability propagation
    - 🙁 License violation
    - 🙁 Supply chain attack
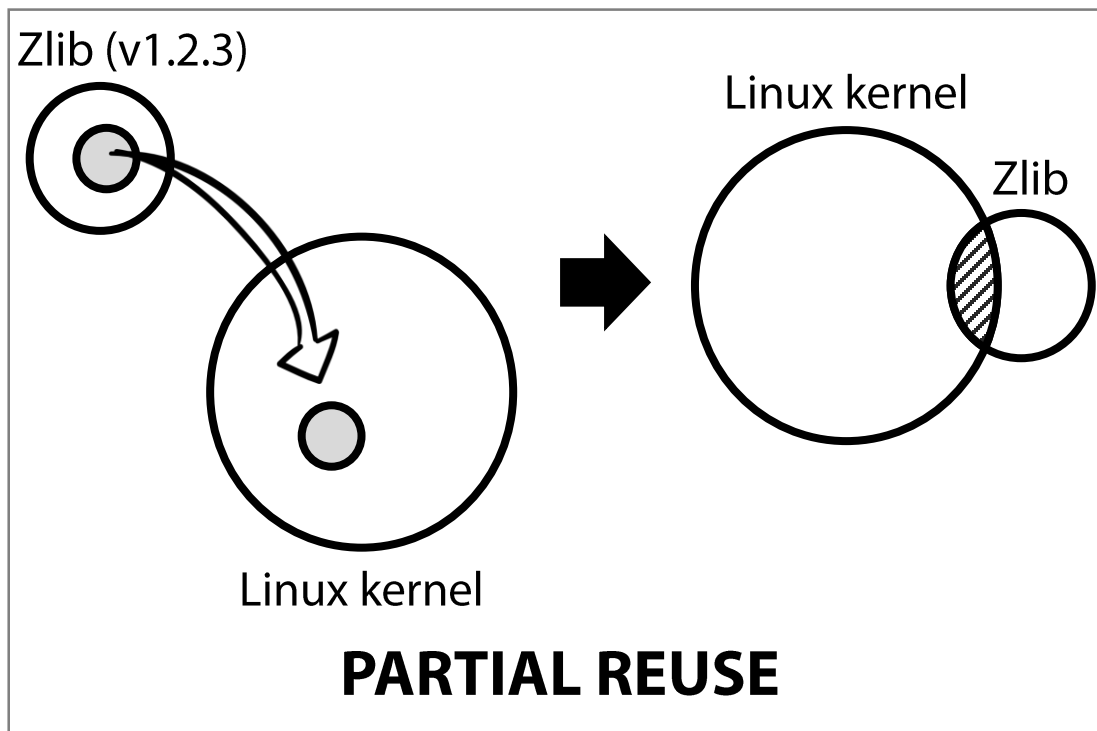
1

# CHALLENGES

- **Previous approaches cannot precisely identify OSS components**

  - **Modified OSS reuse**
    - The cause of <u>false negatives</u> in component identification

  - **Nested OSS components**
    - The cause of <u>false positives</u> in component identification

# CHALLENGES

- **Modified OSS reuse**
  - Modified reuse patterns
    - <u>Partial</u> reuse, <u>structure-changed</u> reuse, <u>code-changed</u> reuse



Zlib (v1.2.3)

Linux kernel

Zlib
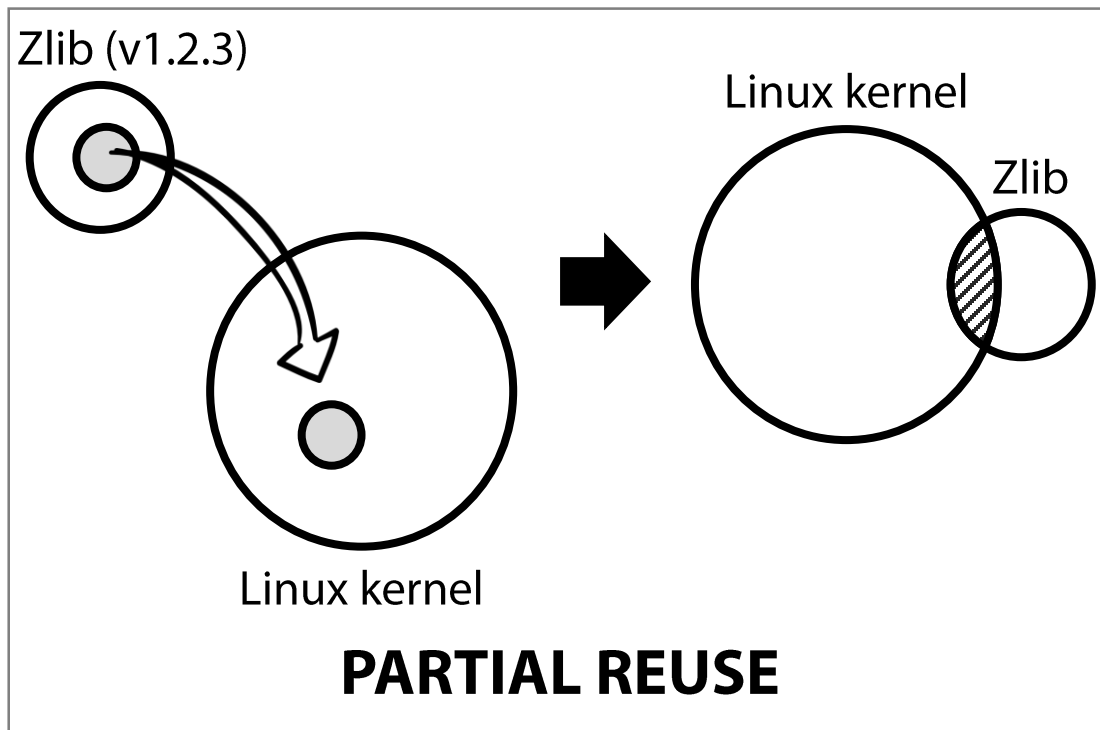
Linux kernel

**PARTIAL REUSE**

```
/* inflate.c -- zlib decompression
 * Copyright (C) 1995-2005 Mark Adler
 * For conditions of distribution and use, see copyright notice in zlib.h
 *
 * Based on zlib 1.2.3 but modified for the Linux Kernel by
```

# CHALLENGES

- **Modified OSS reuse**
  - Modified reuse patterns
    - <u>Partial</u> reuse, <u>structure-changed</u> reuse, <u>code-changed</u> reuse
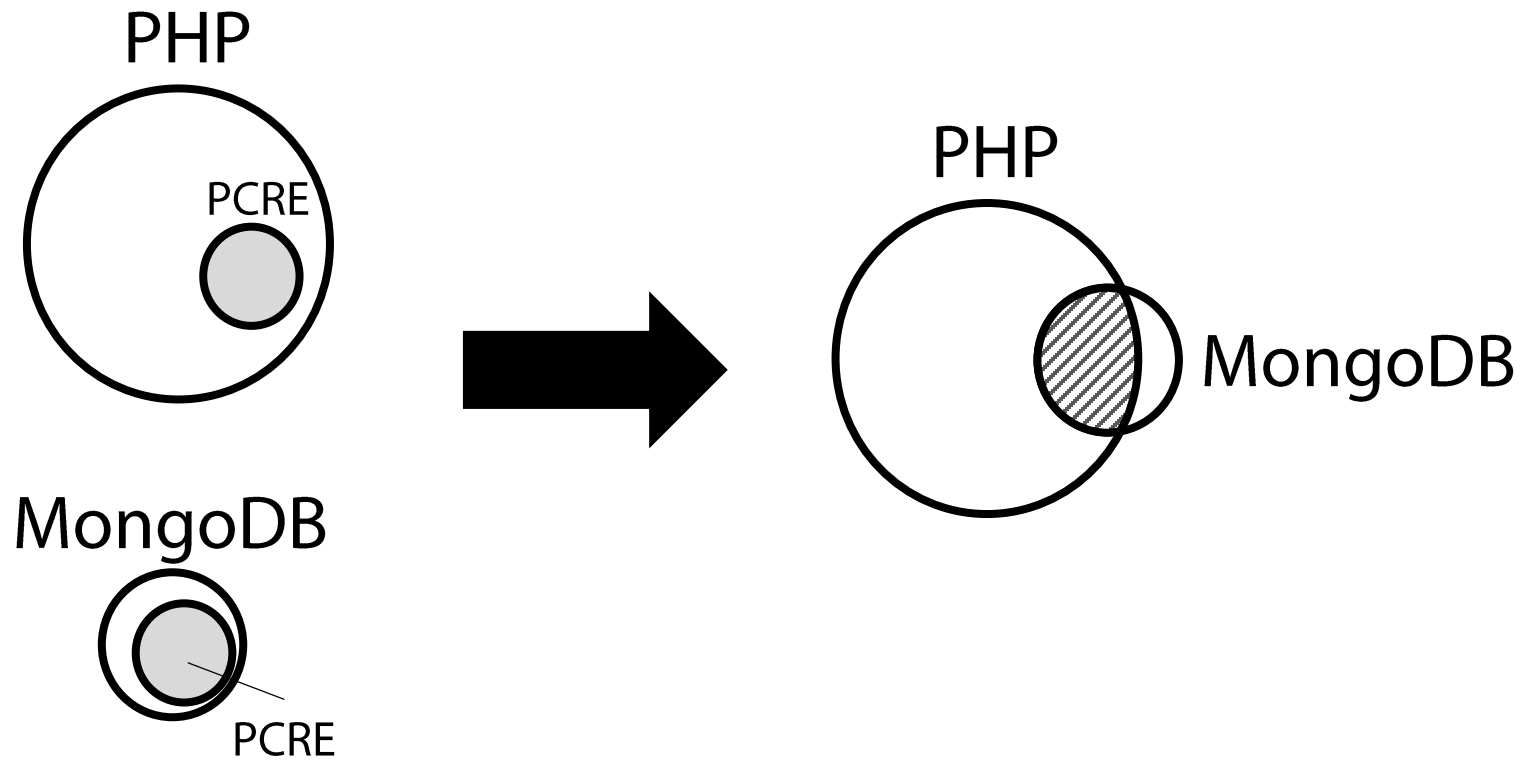


Zlib (v1.2.3)

Linux kernel

Zlib

Linux kernel

**PARTIAL REUSE**

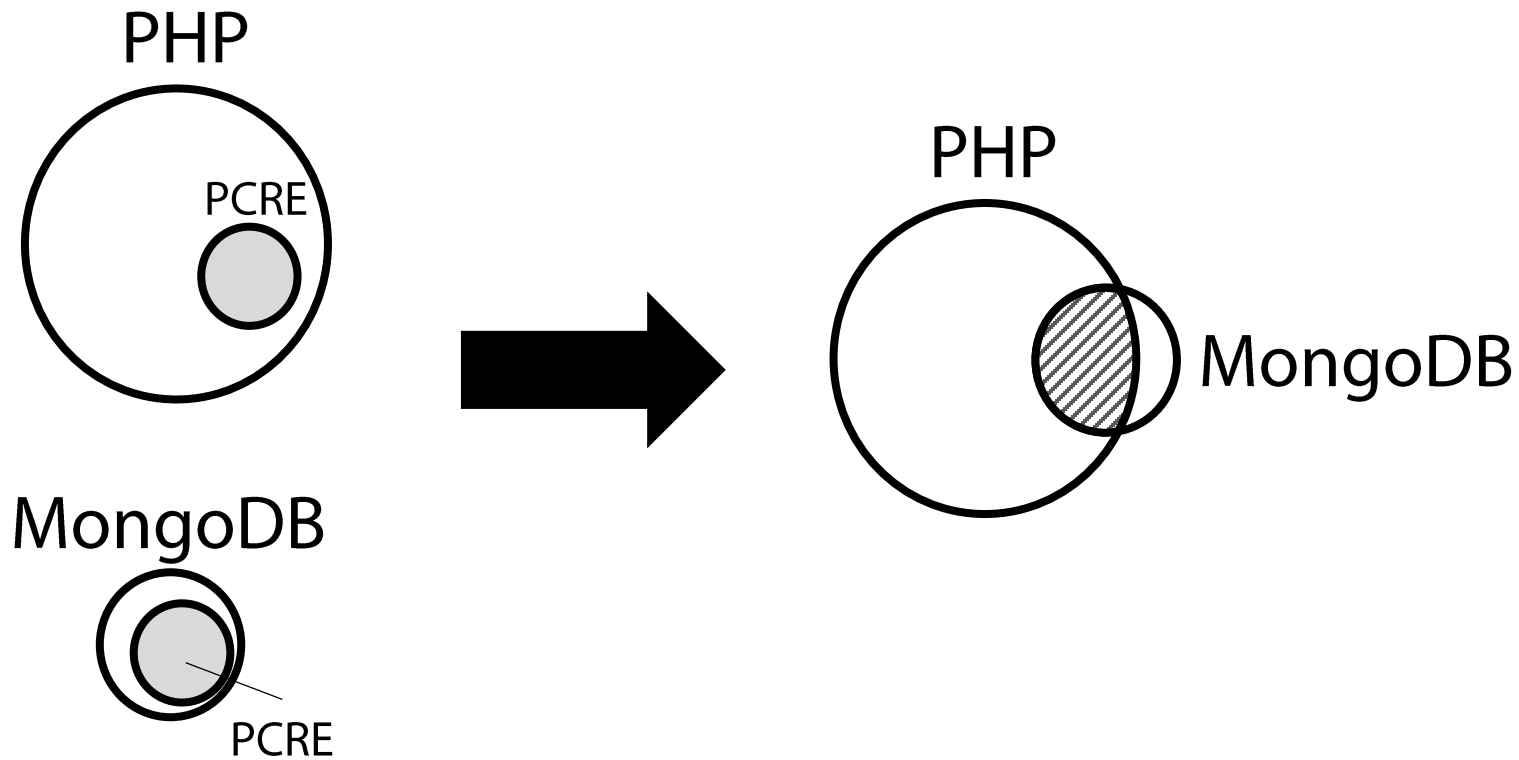Simple threshold-based approach

Many
**false negatives**

3

# CHALLENGES

- **Nested components**

# CHALLENGES

- **Nested components**
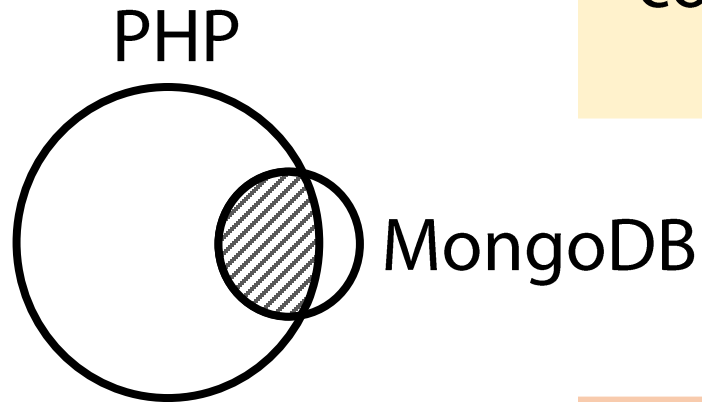


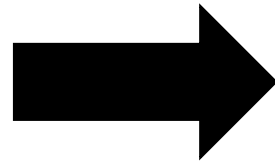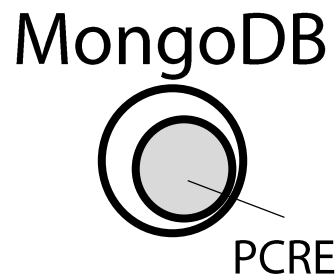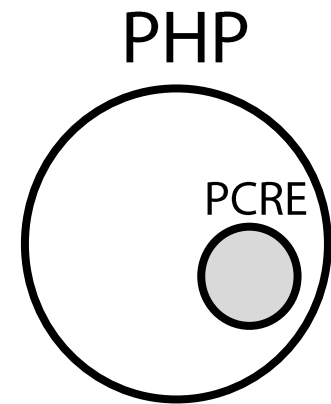PHP

PCRE

MongoDB

PCRE

PHP

MongoDB
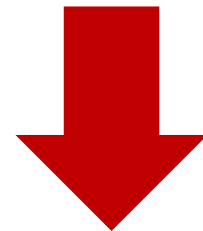
**Correct answers**
- PHP reuses PCRE
- MongoDB reuses PCRE

**Wrong answers**
- MongoDB reuses PHP
- PHP reuses MongoDB

# CHALLENGES

- **Nested components**



PHP

PCRE

MongoDB

PCRE

PHP

MongoDB

Existing software composition analysis approaches

Many
**false positives**

# CENTRIS

- **CENTRI**fuge for **S**oftware

  - The first approach to precisely and scalably identify **modified** OSS components
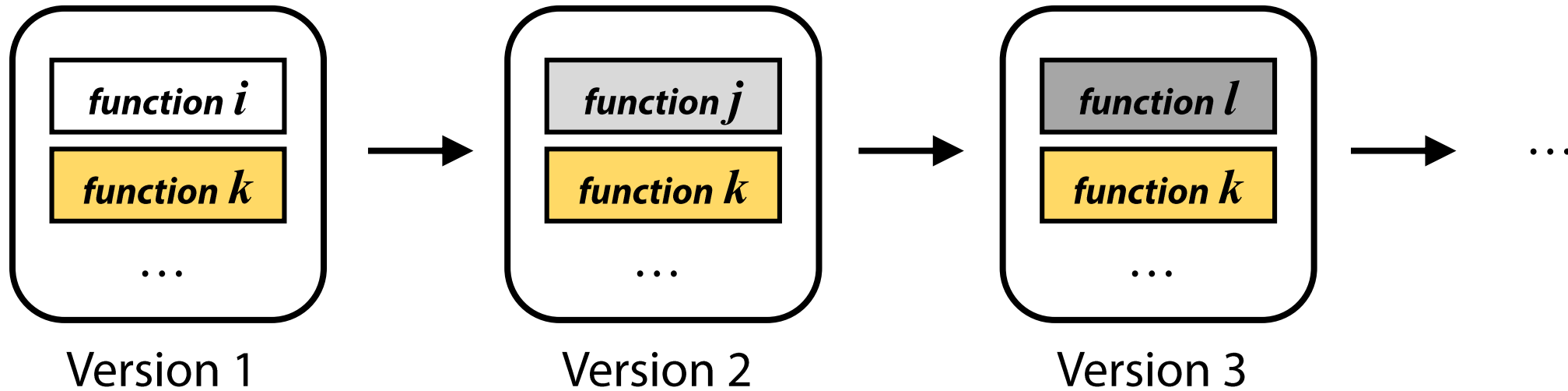
  - Key techniques

    - **S1. Redundancy elimination**

      - For *high scalability*

    - **S2. Code segmentation**
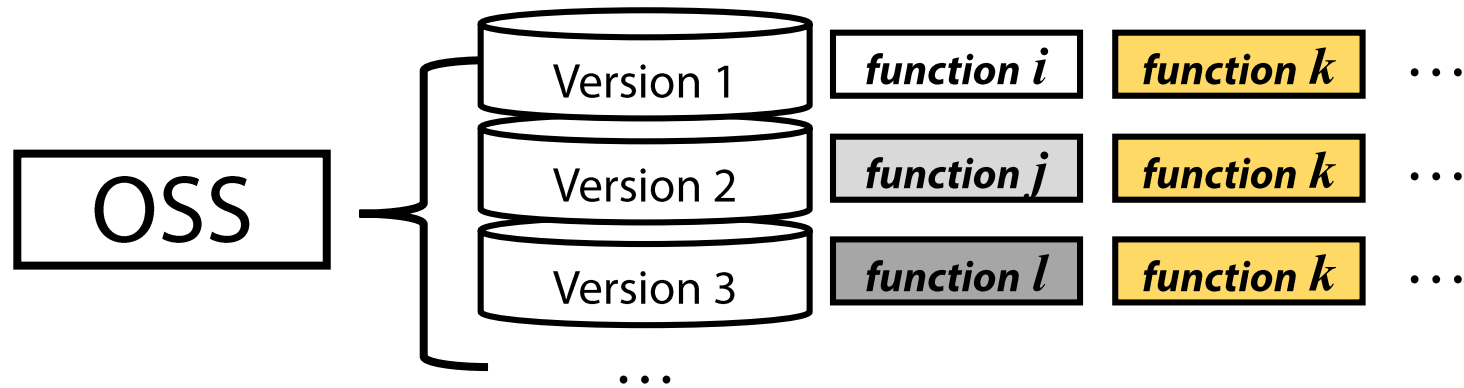
      - For *high accuracy*

5

# S1. Redundancy elimination



Version 1 → Version 2 → Version 3 → …
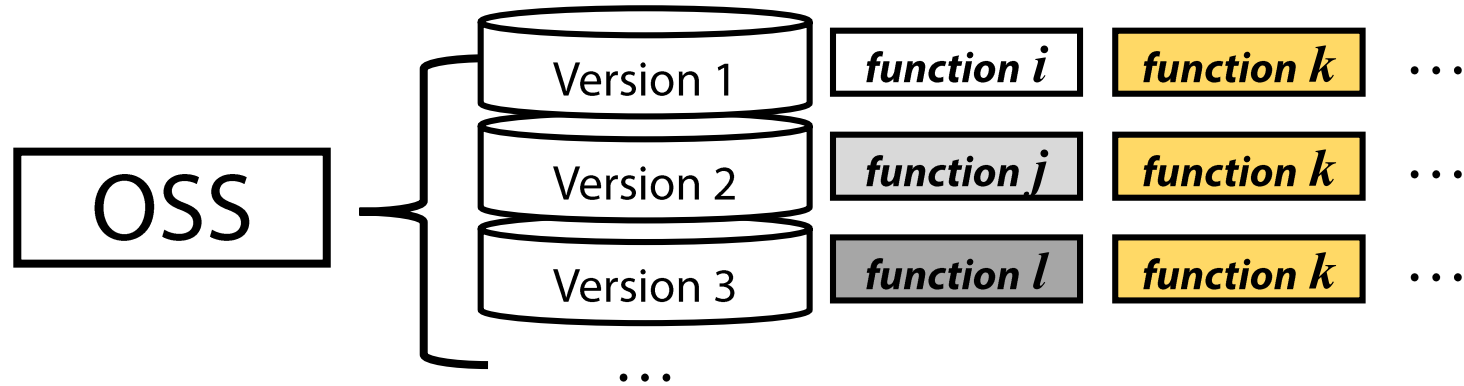
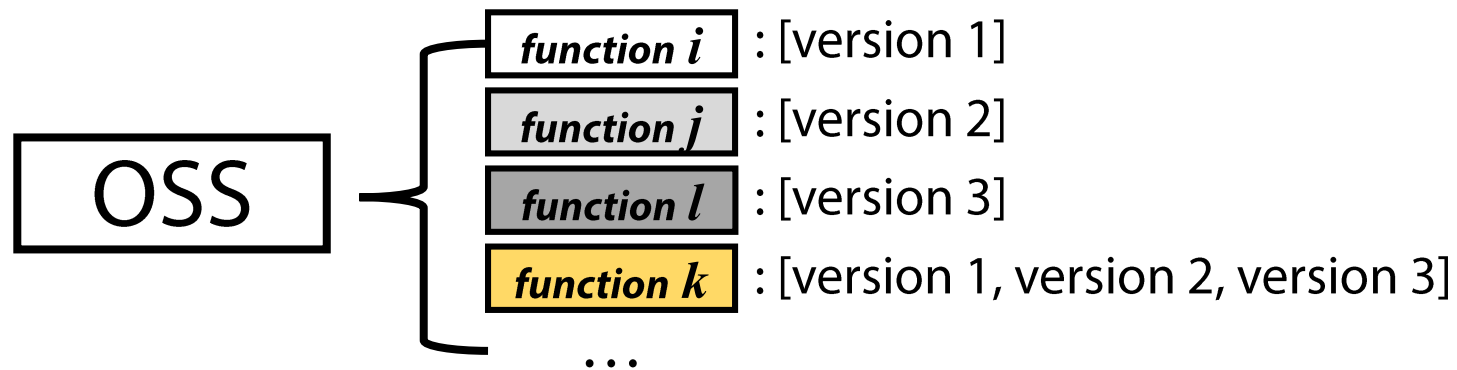Version update in an OSS

# S1. Redundancy elimination



**A naively generated OSS signature**

| function $k$ | : compared 3+ times |

# S1. Redundancy elimination



**A naively generated OSS signature**



**A redundancy eliminated signature for an OSS**

# S1. Redundancy elimination



**A naively generated OSS signature**

**A redundancy eliminated signature for an OSS**

# S2. Code segmentation



OSS

**Application code**

**Borrowed code**

**The unique part of the software**
- Non-reused code parts
- Self-developed code

**Non-unique part of the software**
- Reused code parts
- **Cause of false alarms**

# S2. Code segmentation



Code segmentation

S₁    S₂    S₃

# S2. Code segmentation

- **How to segment an OSS?**

# S2. Code segmentation

- **Detecting functions belonging to the <u>borrowed code part</u> of S**

S          L

Birth date of $f$ in S

OSS (S)

OSS (L)

Birth date of $f$ in L

$f \subset$ Borrowed code part of S

- **Detecting functions belonging to the <u>borrowed code part</u> of S**



$$G = \{f \mid \big(f \in (S \cap L)\big) \wedge (birth(f, L) \leq birth\ (f, S))\}$$

# S2. Code segmentation

1) **Measure similarity between S and L**

$$\phi(\mathrm{S}, \mathrm{L}) = \frac{|G|}{|L|}$$

# S2. Code segmentation

1) **Measure similarity between S and L**

$$\phi(S, L) = \frac{|G|}{|L|}$$

2) **Check whether G is included in the borrowed code part of S**

If $\phi \geq \theta$ then :



or

# S2. Code segmentation

1) **Measure similarity between S and L**

$$\phi(S, L) = \frac{|G|}{|L|}$$

2) **Check whether G is included in the borrowed code part of S**

If $\phi \geq \theta$ then :



or

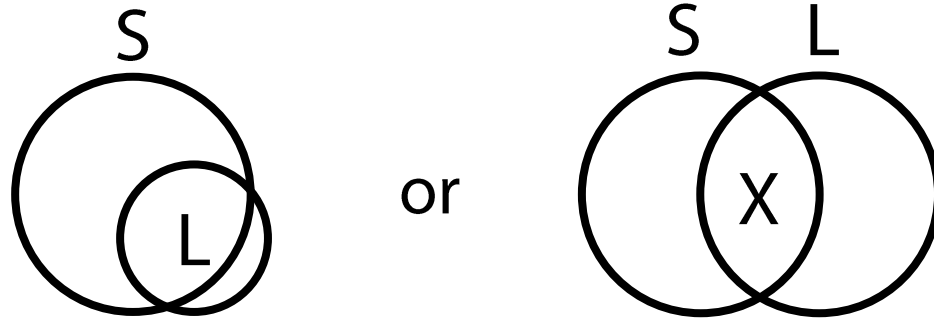3) **Remove G from S**
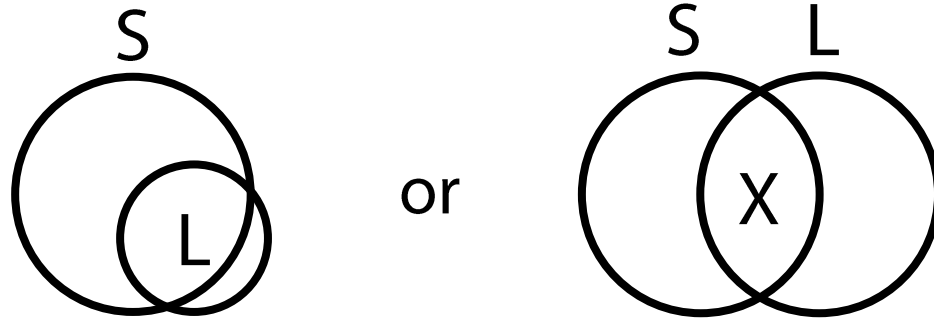
$$S = (S \setminus G)$$

# S2. Code segmentation

1) **Measure similarity between S and L**

$$\phi(S, L) = \frac{|G|}{|L|}$$

2) **Check whether G is included in the borrowed code part of S**

If $\phi \geq \theta$ then :

S

or

S    L

X

3) **Remove G from S**

$S = (S \setminus G)$

**Repeat this process for all OSS in the component DB**

=> Only <u>the application code of S</u> remains

# Component identification in the target software

- **Comparing T with the application code part of the collected OSS**



$$\Phi(\mathrm{T},S) = \frac{|T \cap S_A|}{|S_A|}$$

=> if $\mathbf{\Phi(T, S) \geq \theta}$, then S is the component of T

# EVALUATION

- **Dataset**
  - **Popular C/C++ OSS projects from GitHub (April, 2020)**
    - #Stars >= 100
    - A total of 10,241 projects, 229,326 versions, and 80 billion lines of code (LoC)

- **Parameter**
  - $\theta = 0.1$

14

# Evaluation

1) **Accuracy**
- Cross-comparison experiments (10,241 vs 10,241)
- **91% precision and 94% recall**
  - Modified components account for 95% of the detected components!

# EVALUATION

**1) Accuracy**

- Cross-comparison experiments (10,241 vs 10,241)
- **91% precision and 94% recall**
    - Modified components account for 95% of the detected components!

**2) Scalability**



□—□ SourcererCC　　○┄┄○ CENTRIS (first exp.)　　●—● CENTRIS ($n^{th}$ exp.)

Time (in hours): 800, 600, 400, 200, 0

The limitation of SourcererCC due to memory error →

98 hours (first exp.)
↓
A minute ($n^{th}$ exp.)

Dataset (LoC): 1M, 10M, 100M, 1B, 5B
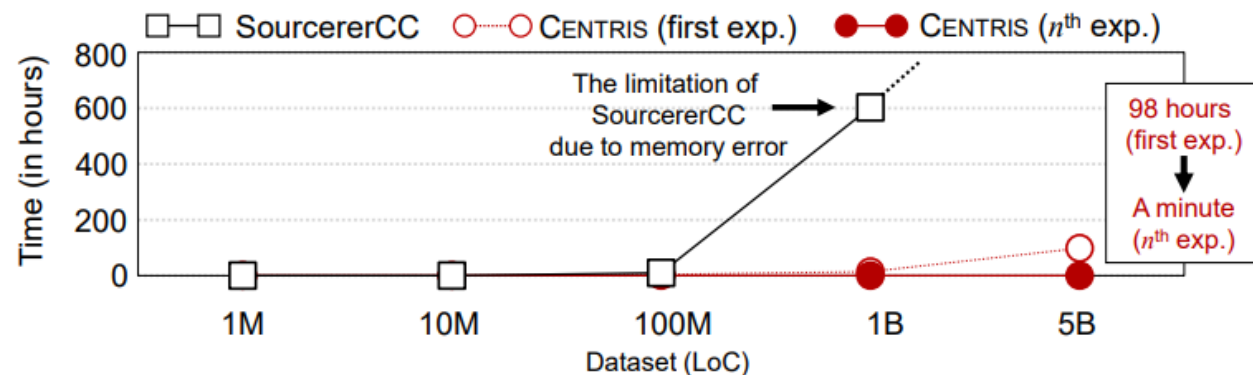
15

# EVALUATION

**1) Accuracy**

- Cross-comparison experiments (10,241 vs 10,241)
- **91% precision and 94% recall**
  - Modified components account for 95% of the detected components!

**2) Scalability**



**3) Identification speed**

- Takes ≤ **1 min** to identify components in the 1 M LoC target software
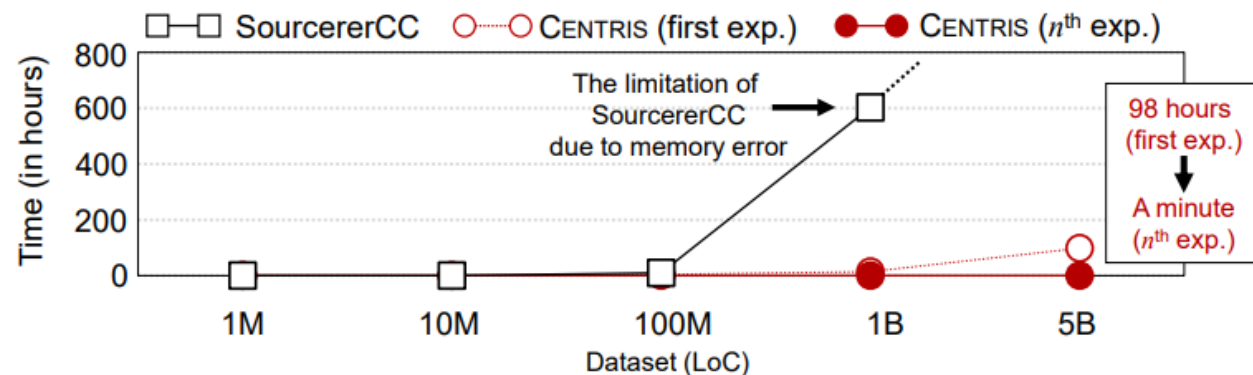
15

# EVALUATION

## 1) Accuracy

- Cross-comparison experiments (10,241 vs 10,241)
- **91% precision and 94% recall**
  - Modified components account for 95% of the detected components!
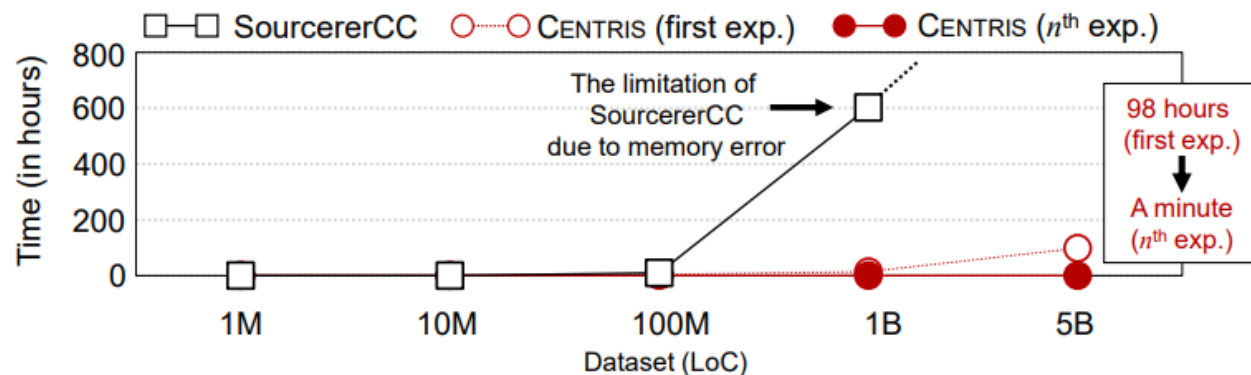
## 2) Scalability



## 3) Identification speed

- Takes ≤ **1 min** to identify components in the 1 M LoC target software

## 4) vs. DejaVu (OOPSLA 2017)

- Code-duplication detection tool
- Using four target software programs
- DejaVu showed only 10% precision

|  | DejaVu | CENTRIS |
|---|---|---|
| Precision | 10% | 95% |
| Recall | 40% | 100% |

# CONCLUSION

- **95% of detected components were reused with modification**
    - Modified components, not likely to be identified, have more chances to pose security threats
    - Management for supply chains considering modified components is required


- **CENTRIS can be the first step towards addressing problems arising from unmanaged OSS components in practice**
    - With the information provided by CENTRIS, developers can mitigate security threats
        - e.g., they can update old-and-vulnerable components

16

## Q&A

Thank you for your attention!

- CENTRIS repository (https://github.com/wooseunghoon/Centris-public)
- CENTRIS at IoTcube (https://iotcube.net/Centris)

## CONTACT

- Seunghoon Woo (seunghoonwoo@korea.ac.kr, https://wooseunghoon.github.io)
- Computer & Communication Security Lab (https://ccs.korea.ac.kr)
- Center for Software Security and Assurance (https://cssa.korea.ac.kr)