# UNIVERSITÉ DE LORRAINE





(00001011) Laboratoire lorrain de recherche (00001011) en informatique et ses applications

# Scalable Program Clone Search through Spectral Analysis

Tristan Benoit benoit.tristan.info@gmail.com Université de Lorraine, CNRS, LORIA Nancy, France Jean-Yves Marion jean-yves.marion@loria.fr Université de Lorraine, CNRS, LORIA Nancy, France Sébastien Bardin sebastien.bardin@cea.fr CEA LIST, Université Paris-Saclay Saclay, France

### 11000010110

ESEC/FSE 2023 Sun 3 - Sat 9 December 2023 San Francisco, California, United States



### **Program Clones**



GCC 4.8.5 -00

GCC 9.3.0 -00

same source code but different toolchains / optimizations / architecture

same source code but obfuscation

ESEC/FSE 2023

(slightly) different version of the source code

### **Applications**

#### **Studying Malware**

Genealogy

01000010101000010

100001011

- Clustering
- Detection

#### Detecting

- Software Theft
- Plagiarism

#### Identifying libraries

- Software Engineering
- Cybersecurity (1day)

Malware – supply chain – firmware – closed programs

ESEC/FSE 2023



### Tsunami for ARM 32 genealogy.

« The Tangled Genealogy of IoT Malware » Cozzi, Vervier, Dell'Amico, Shen, Balzarotti. ACSAC 2020



#### 1) Preprocessing

We extract features from target program P.

#### 2) Similarity Checks

For each program Q in repository R, compute the similarity index between P and Q.

#### 3) Decision

01101100

00001011

We **pick the program** Qm with the **highest similarity index**. The clone search is a success if Qm is indeed a clone of P.

ESEC/FSE 2023 Scalable Program Clone Search through Spectral Analysis

# **A Blooming Field**

Approach	Year	Venue	Input	Approach	Year	Venue	Input
TRACY [1]	2014	PLDI	F	<b>BINSEQUENCE</b> [8]	2017	ASIACCS	F
BINCLONE [65]	2014	SERE	I*	XMATCH [9]	2017	ASIACCS	F
RMKNHLLP2014 [66]	2014	DIMVA	F*	CACOMPARE [74]	2017	ICPC	F
CXZ2014 [21]	2014	TDSC	P	SPAIN [30]	2017	ICSE	P
BLEX [67]	2014	USENIX Sec	F	BINSIGN [75]	2017	IFIP SEC	F
CoP [33], [68]	2014	ESEC/FSE	P	GITZ [10]	2017	PLDI	F
TEDEM [2]	2014	ACSAC	B*	BINSHAPE [76]	2017	DIMVA	F
SIGMA [69]	2015	DFRWS	F	BINSIM [77]	2017	USENIX Sec	T
MXW2015 [24]	2015	IFIP SEC	P	KS2017 [31]	2017	ASE	T
MULTI-MH [3]	2015	S&P	B*	IMF-SIM [78]	2017	ASE	F
QSM2015 [70]	2015	SANER	F	GEMINI [12]	2017	CCS	F
DISCOVRE [4]	2016	NDSS	F	Fossil [79]	2018	TOPS	F
MOCKINGBIRD [29]	2016	SANER	F	FIRMUP [13]	2018	ASPLOS	F
Esh [5]	2016	PLDI	F	BINARM [14]	2018	DIMVA	F
Трм [71]	2016	TrustCom	P	$\alpha \text{DIFF}$ [15]	2018	ASE	P
BinDnn [72]	2016	SecureComm	F	VULSEEKER [11]	2018	ASE	F
GENIUS [6]	2016	CCS	F	RLZ2019 [80]	2019	BAR	B
BinGo [7]	2016	FSE	F	INNEREYE [81]	2019	NDSS	B*
KLKI2016 [18]	2016	JSCOMPUT	P	ASM2VEC [82]	2019	S&P	F
KAM1N0 [73]	2016	SIGKDD	B*	SAFE [83]	2019	DIMVA	F

#### I:Instruction B:Basic Block F:Function P:Program

« A survey of binary code similarity » Irfan Ul Haq and Juan Caballero. ACM Computing Surveys 2021



ESEC/FSE 2023

010000101000010

101101 PP881 der pag of 1681 199

11011000

1011010

01101100

01101001

Scalable Program Clone Search through Spectral Analysis

# **A Blooming Field**

Precision

881der pag

0110110

110110

Robustness (cross-compilers, cross-archi., light obfuscation)
 Speed/Scale (size of code, size of repository)

CoP [33], [68]	2014	ESEC/FSE	Р	GITZ [10]	2017	PLDI	F
TEDEM [2]	2014	ACSAC	B*	BINSHAPE [76]	2017	DIMVA	F
SIGMA [69]	2015	DFRWS	F	BINSIM [77]	2017	USENIX Sec	T
MXW2015 [24]	2015	IFIP SEC	Р	KS2017 [31]	2017	ASE	T
MULTI-MH [3]	2015	S&P	B*	<b>IMF-SIM</b> [78]	2017	ASE	F
QSM2015 [70]	2015	SANER	F	GEMINI [12]	2017	CCS	F
DISCOVRE [4]	2016	NDSS	F	Fossil [79]	2018	TOPS	F
MOCKINGBIRD [29]	2016	SANER	F	FIRMUP [13]	2018	ASPLOS	F
Esh [5]	2016	PLDI	F	BINARM [14]	2018	DIMVA	F
ТРМ [71]	2016	TrustCom	Р	$\alpha \text{DIFF}$ [15]	2018	ASE	P
BINDNN [72]	2016	SecureComm	F	VULSEEKER [11]	2018	ASE	F
GENIUS [6]	2016	CCS	F	RLZ2019 [80]	2019	BAR	B
BINGO [7]	2016	FSE	F	INNEREYE [81]	2019	NDSS	<b>B</b> *
KLKI2016 [18]	2016	<b>JSCOMPUT</b>	Р	ASM2VEC [82]	2019	S&P	F
KAM1N0 [73]	2016	SIGKDD	B*	SAFE [83]	2019	DIMVA	F

#### I:Instruction B:Basic Block F:Function P:Program

« A survey of binary code similarity » Irfan Ul Haq and Juan Caballero. ACM Computing Surveys 2021



# **A Blooming Field**

Precision

PPBBIder page

0110110

110110

- Robustness (cross-compilers, cross-archi., light obfuscation)
  Speed/Scale (cize of code, cize of repeatery)
- Speed/Scale (size of code, size of repository)

ESEC/FSE 2023



### I:Instruction B:Basic Block F:Function P:Program

« A survey of binary code similarity » Irfan Ul Haq and Juan Caballero. ACM Computing Surveys 2021

### A case for program-level clone search

### Claim:

### Clone search is essentially at the level of programs (not functions)

- Malware studies, software theft, library identification
- A single function is not the right focus

ESEC/FSE 2023

#### Function-level methods do not scale to program-level

- Gemini: 17h for 20 clone searches on 1420 libraries
- SAFE: 160h for the same 20 clone searches

#### **GED and Matching methods do not scale with code size**

- SMIT: 43h to compute a single similarity index between Geany and the cp command
- DeepBinDiff: 10min for basic block matching on small binaries from core utils



# A case for program-level clone search

### **Goal: a program-level clone search approach**

- Scalable
- Reasonably precise and robust (cross archi / compilers, etc.)
- Do not rely on symbols and strings

ESEC/FSE 2023

### Contributions

- PSS (Program Spectral Similarity), a similarity metric between programs that is fast and accurate while being very robust even against (light) obfuscations
  - key insight: spectral graph analysis
  - Dedicated adaptations and optimization (PSSO)

No learning

- An evaluation covering 14 methods and over 200,000 programs from Linux, Windows and IoT malware
- Insights about how classes of methods perform

Available at: <a href="https://doi.org/10.5281/zenodo.8289599">https://doi.org/10.5281/zenodo.8289599</a>

ESEC/FSE 2023

id PPBArder page

01101100

oric



Scalable Program Clone Search through Spectral Analysis

Try it!

### PSS

### Ideas

o o PPBBIder pag

- Performing quick comparisons with the spectral distance between graph representations of programs.
- Reducing the one-time computation of the spectrum by using only the spectrum of the call graph.
- Incorporating information about each function CFG.





**Definition** An eigenvalue  $\lambda$  and an eigenvector  $\vec{u}$  is a solution to the equation :  $(L - \lambda I) \vec{u} = \vec{0}$ . The spectrum is the set  $\Lambda$  of eigenvalues  $\{\lambda_1(G), \lambda_2(G), \ldots, \lambda_{|G|}(G)\}$  where  $\lambda_1(G) \ge \lambda_2(G) \ge \ldots \ge \lambda_{|G|}(G)$  and where |G| is the number of vertices of G.

ESEC/FSE 2023

### PSS

**Definition** The spectral distance [54] between two graphs  $G_1$ ,  $G_2$  of same order n is  $\sum_{i=1}^{N} |\lambda_i(G_1) - \lambda_i(G_2)|.$ The generalized spectral distance is defined by  $: \sum_{i=1}^{\min(|G_1|,|G_2|)} |\lambda_i(G_1) - \lambda_i(G_2)|.$ 

**Theorem** Let G' be a copy of a graph G except for a removed vertex of degree r. For all i such that  $1 \le i \le n-1$ ,  $\lambda_i(G) \ge \lambda_i(G') \ge \lambda_{i+r}(G)$ .

**Corollary** The generalized spectral distance between G and G' is equal to a 2r.

ESEC/FSE 2023

101101010999940010000119814000

0110110

0110110

The Lanczos algorithm calculate the full spectrum in  $O(dn^2)$  where n is the number of vertices, and d is the average degree of the graph.



- 1) From an undirected version of the **call graph**, we compute its spectrum **v** and **normalize it**.
- For each function, we compute its number of edges. We gather it as a vector w, sort it in descending order, and normalize it.

Similarity index between programs in O(n):

$$simCG(P_0, P_1) := \sqrt{2} - \sqrt{\sum_{i=0}^{\min(|\vec{w_0}|, |\vec{w_1}|)} (v_{0,i} - v_{1,i})^2}$$
$$simCFG(P_0, P_1) := \sqrt{2} - \sqrt{\sum_{i=0}^{\min(|\vec{w_0}|, |\vec{w_1}|)} (w_{0,i} - w_{1,i})^2}$$

$$PSS(P_0, P_1) := \frac{simCG(P_0, P_1) + simCFG(P_0, P_1)}{2\sqrt{2}}$$

ESEC/FSE 2023

o 14 PP881 der pag

0110110

# **The PSSO optimization**

We propose PSSO, an optimized version of PSS. Instead of computing the whole spectrum, we compute only the **100 greatest eigenvalues** of the call graph.

Lanczos algorithm can compute the k greatest eigenvalues in O(k.d.n) where n is the number of vertices and d the average degree of the graph.

2891derpa

### **Theoretical complexities**

011010010110000101

011000110

101100

00101101000

01101100

01101111 01110010 01101001

00001011

Method	Class	Similarity Check <sup>†</sup>	Preprocessing <sup>‡</sup>
SMIT [45]	GED	$O(n^4)$	O(dn)
CGC [104]	Matching	$O(n^4)$	O(dn)
MutantX-S [44]	N-gram	<i>O</i> (1)	O(i)
Asm2Vec [24]	Function ML	$O(n^2)$	O(n)
Gemini [105]	Function ML	$O(n^2)$	O(n)
SAFE [71]	Function ML	$O(n^2)$	O(n)
$\alpha$ Diff [67]	Function ML	$O(n^2)$	O(n)
LibDX [96]	Literal	O(s)	O(s)
LibDB [97]	Literal and Function ML	$O(n^2 + s)$	O(s+n)
DeepBinDiff [26]	Matching and ML	$O(n^3m^3)$	None
PSS	Spectral	O(n)	$O(dn^2)$
PSSO	Spectral	O(n)	O(dn)

n: # functions, i: # instructions , s: # string literals

d : average # call per functions, m : # basic blocks per function

† Between two programs

‡ Once for the whole clone search

ESEC/FSE 2023

#### Scalable Program Clone Search through Spectral Analysis

# **Preliminary Study**

Total time for all clone searches on the Basic dataset

$\mathbf{B}_{size}$	$\leq 1h30m$
$\mathbf{D}_{size}$	$\leq 1h30m$
Shape	$\leq 1 h30 m$
ASCG	$\leq 1h30m$
MutantX-S	$\leq 1 h30 m$
PSS	$\leq 1h30m$
PSSO	$\leq 1h30m$
LibDX	$\leq 1h30m$
StringSet	$\leq 1h30m$
FunctionSet	$\leq 1h30m$

100 10 PPBBH der pag bis

0110110

110110

ASCFG	42h
GED-0	81h
GED-L	46h
SMIT	3634h
CGC	171h
Asm2vec $\dagger$	141h
Gemini †	102h
SAFE †	655h
$\alpha \text{Diff}$ †	642h
LibDB †	16h

Learning time not included

10 Methods selected for comparison

### Focus of the talk: the 7 methods that do not rely on strings

ESEC/FSE 2023

### BinKiT

From 51 GNU packages, 235 source codes have been extracted. They have been compiled with 416 different toolchains for a total of 97,760 programs.

#### It cover:

19 PPBBIder page

0110110

- 8 architectures (arm, x86, mips, and mipseb, for 32/64 bits)
- 9 compilers (5 versions of GCC and 4 versions of Clang)
- 4 optimization levels from -00 to -03
- 4 obfuscations from Obfuscator-LLVM (SUB, BCF, FLA, all 3)

*« Revisiting Binary Code Similarity Analysis using Interpretable Feature Engineering and Lessons Learned ».* 

*Dongkwan Kim, Eunsoo Kim, Sang Kil Cha, Sooel Son, Yongdae Kim IEEE Transactions on Software Engineering. 2022* 

ESEC/FSE 2023

#### **IoT Malware**

19,959 IoT malwares from MalwareBazaar (2020-2022) 3 clones family:

Arm: 41.6%

8 major

architectures

- 12,357 mirai
- 5842 gafgyt

uperH: 8.0

MIPS: 21.1%

x86: 7.2%

SPARC : 5.6%

• 1760 tsunami

ARC: 0.5% IBM System: 0.0%

Motorol 68040:

ESEC/FSE 2023

### Windows

84,992 programs for Windows Clones share:

1) filename 2) target platform

49,443 programs have a clone.



1010010110000

1199 PP881 der pag at 1841

### Clone search average runtime ( preprocessing )

Dataset	Basic	BinKiT	IoT	Windows
# Programs	1K	96K	20K	85K
B <sub>size</sub>	< 0.01	0.11	0.14	0.63
D <sub>size</sub>	< 0.01	0.11	0.14	0.63
Shape	0.02	0.05	0.06	0.31
ASCG	1.42 (1.42)	0.37 (0.21)	0.25 (0.06)	17.68 (16.60)
MutantX-S	< 0.01	0.57	0.64	3.00
PSS	1.41 (1.41)	0.49 (0.21)	0.39 (0.05)	19.17 (16.95)
PSSO	0.27 (0.27)	0.30 (0.04)	0.44 (0.14)	2.29 (0.39)

PSS is fast, except on Windows. PSSO solves this issue. MutantX-S is also very fast.

ESEC/FSE 2023

1d PPBBHder and by

0110110

0110110

#### Average success

Dataset	Basic	BinKit	IoT	Windows
B <sub>size</sub>	0.17	0.166	0.819	0.196
$D_{size}$	0.16	0.062	0.787	0.445
Shape	0.19	0.297	0.818	0.389
ASCG	0.24	0.554	0.759	0.444
MutantX-S	0.38	0.354	0.870	0.472
PSS	0.38	0.619	0.863	0.475
PSS <sub>O</sub>	0.38	0.619	0.862	0.466

PSS/PSSO are as accurate as MutantX-S, except on BinKiT on which they are better.

ESEC/FSE 2023

o Paged on and

#### Average success on BinKiT

Category		Optimization level					Cross-compiler			Cross-architecture				v	s. Obfu	scatior	1†
	O0	O0	O0	01	O1	O2	gcc-4	clang-4	clang	arm	arm	mips	32				
VS.	01	O2	O3	O2	O3	O3	gcc-8	clang-7	gcc	mips	x86	x86	64	bcf	fla	sub	all
B <sub>size</sub>	0.04	0.04	0.07	0.19	0.11	0.21	0.11	0.45	0.07	0.03	0.10	0.04	0.04	0.04	0.01	0.08	0.01
$D_{size}$	0.03	0.03	0.03	0.06	0.05	0.07	0.07	0.09	0.04	0.02	0.05	0.03	0.04	0.02	0.01	0.05	0.01
Shape	0.19	0.07	0.06	0.17	0.11	0.33	0.38	0.65	0.16	0.04	0.16	0.04	0.19	0.25	0.27	0.48	0.23
ASCG	0.40	0.12	0.10	0.43	0.24	0.68	0.78	0.91	0.46	0.08	0.46	0.06	0.59	0.54	0.64	0.78	0.48
MutantX-S	0.04	0.03	0.03	0.43	0.36	0.64	0.67	0.80	0.14	0.02	0.01	0.01	0.06	0.09	0.03	0.54	0.01
PSS	0.54	0.23	0.17	0.59	0.38	0.70	0.79	0.91	0.51	0.39	0.55	0.39	0.66	0.53	0.57	0.82	0.46
PSS <sub>O</sub>	0.53	0.24	0.17	0.60	0.39	0.68	0.78	0.90	0.51	0.44	0.54	0.44	0.66	0.52	0.56	0.82	0.46

PSS and PSSO are very robust, even in cross-compiler/cross architecture scenarios and against obfuscations.

ESEC/FSE 2023

Scalable Program Clone Search through Spectral Analysis

### Conclusion



Method	speed	precision	robust.
ASCG [27]	+	-	+
MutantX-S [33]	+	+	
PSS/PSS <sub>O</sub>	+/++	+	+

+ Insigts about strings based methods in the paper

### Limitations

PBPIder page

0110110

0110000

01101100

11001001

- 1) Hard to deal with changes in function call graphs
- 2) Evaluated mostly on C/C++ programs

ESEC/FSE 2023

### Thank you for your attention!

Scalable Program Clone Search through Spectral Analysis