

Daniel Plohmann, Manuel Blatt, and Daniel Enders | 2023-04-13



MCRIT: The MinHash-based Code Relationship & Investigation Toolkit

Introduction

\$whoami

- Security Researcher @ Fraunhofer FKIE & University of Bonn

- Research Scope:
 - Analysis of malicious software (malware) / reverse engineering / analysis automation

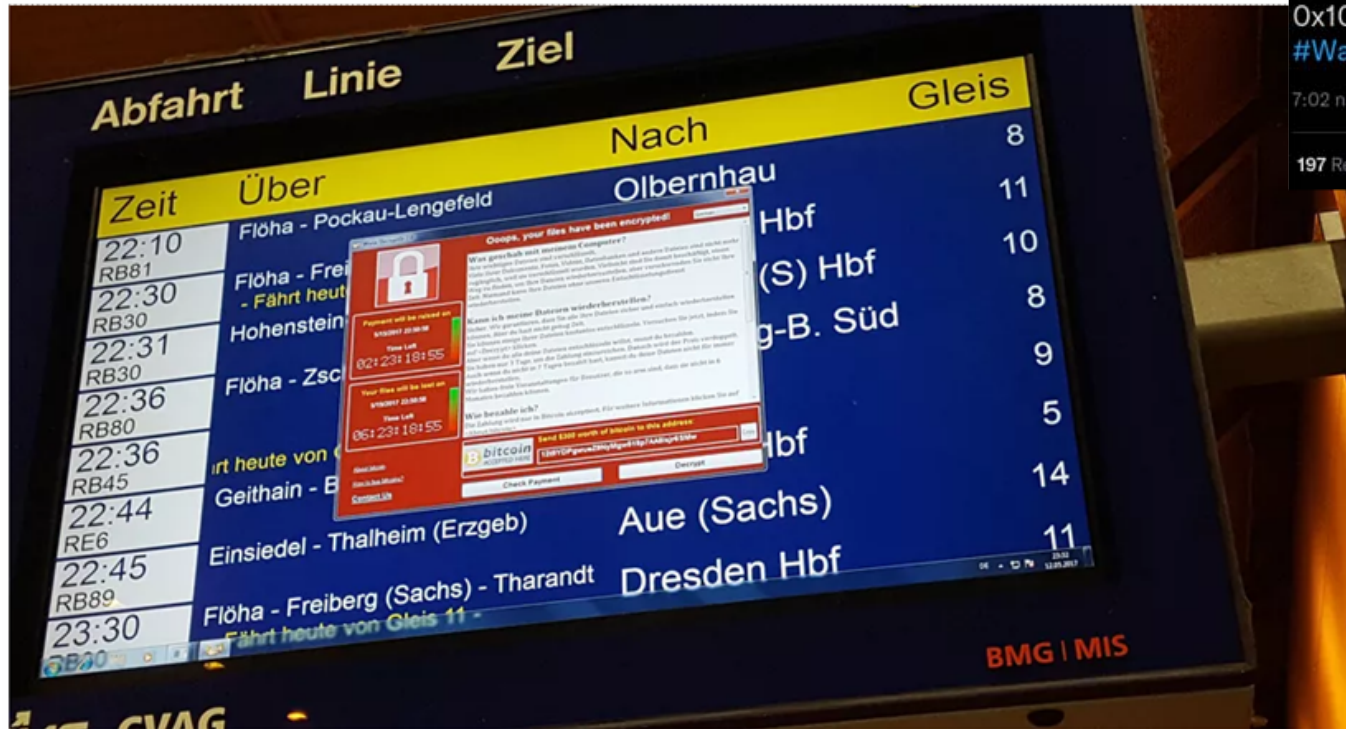
- Past Appearances at Botconf:
 - Laura Guevara, Daniel Plohmann
2014 - Semantic Exploration of Binaries
 - Daniel Plohmann
2015 - DGArchive: A deep dive into domain generating algorithms
 - Daniel Plohmann, Martin Clauß, Steffen Enders, Elmar Padilla
2017 - Malpedia: A Collaborative Effort to Inventorize the Malware Landscape
 - Daniel Plohmann, Steffen Enders, Elmar Padilla
2018 - Code Cartographer's Diary
 - Felix Bilstein, Daniel Plohmann
2019 - YARA-Signator: Automated Generation of Code-based YARA Rules

Outline

- Motivation
- MCRIT: System Overview
 - Methodology
 - Framework
- Use Cases
- Outlook

Motivation

Motivation



Anzeige im Hauptbahnhof in Chemnitz am Freitag: Forderung von Lösegeld Foto: P. Götzel/ dpa



- Infamous WannaCry
- Ransomware attack using wormable exploit (EternalBlue)
 - Attack started on May 12th 2017
 - 230k affected systems in ~8 hours
 - Quickly disrupted due to a lucky registration of killswitch domain
- Impact
 - UK NHS disrupted (£100m damage)
 - Nissan, Renault, Telefonica, FedEx, DB, ...
- Attack Attribution?

[1] https://en.wikipedia.org/wiki/WannaCry_ransomware_attack

[2] <https://twitter.com/neelmehta/status/864164081116225536>

Motivation

[illegible]

Unique across all of Malpedia

[illegible]

 **Neel Mehta**
@neelmehta

9c7c7149387a1c79679a87dd1ba755bc @ 0x402560,
0x40F598
ac21c8ad899727137c4b94458d7aa8d8 @
0x10004ba0, 0x10012AA4
[#WannaCryptAttribution](#)

7:02 nachm. · 15. Mai 2017 · Twitter Web Client

197 Retweets **45** Zitierte Tweets **292** „Gefällt mir“-Angaben

[1] <https://twitter.com/neelmehta/status/864164081116225536>

Motivation

- Code Similarity Analysis
 - High potential to help analysts and accelerate analysis
 - Code identification, library filtering, hunting, label transfer, ...
 - Existing solutions mostly limited to
 - 1:1 comparison
 - Proprietary
- Let's see what we can do. :)

MCRIT

System Overview

MCRIT

Design

- Goal: Analyze code sharing and third-party library usage in malware
 - Create tools to leverage Malpedia binary corpus
 - Don't reinvent the wheel: **reuse** of **proven techniques** as described in literature
- Requirements:
 - Similarity: reliable, interpretable estimate
 - Scalability: (tens of) millions of functions
 - Efficient representation: (significantly) smaller than indexed code
 - No cross-bitness or cross-architecture (*Malpedia at the time was 95% 32bit Intel*)

Code Recovery and Similarity Analysis

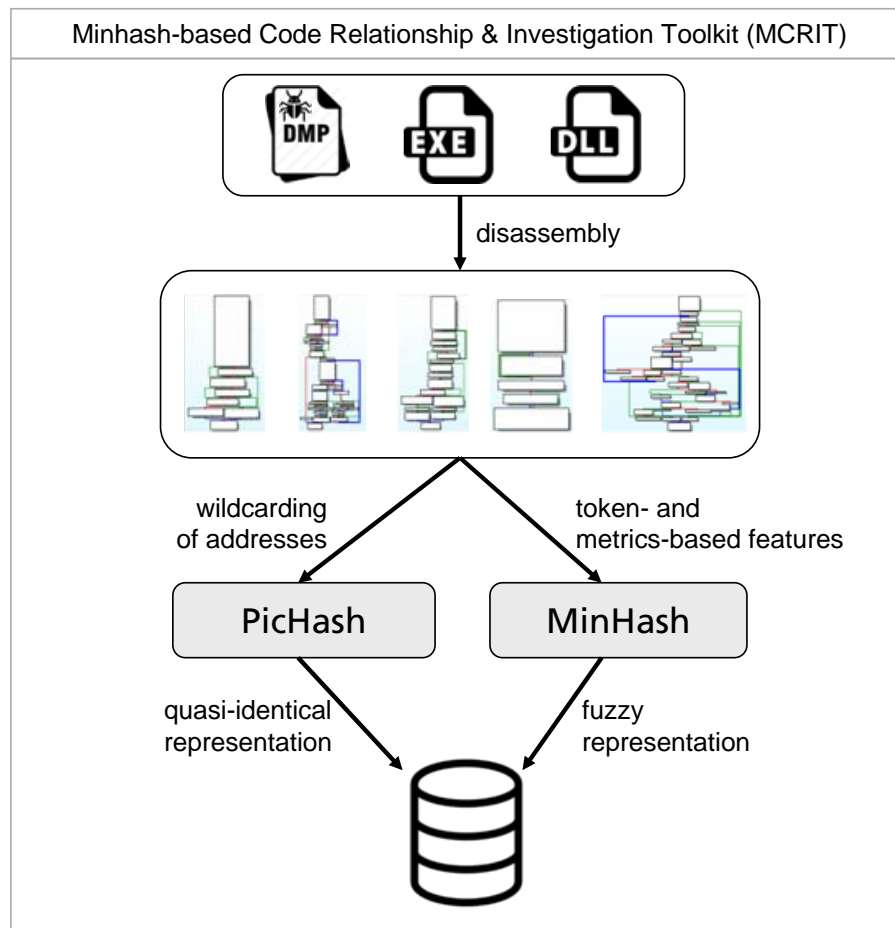
MCRIT: Approach

■ Initial Observations

- Haq et al. [1] survey:
50+ works on code similarity since 2010
- Only 17 with 10+ malware samples
- Only 1 analyzed FOSS usage in malware (Alrabaee et al. [2])

■ MCRIT

- Combines quasi-identical and fuzzy code representation
- Block & Function-level similarity
- Efficient 1:n matching via
 - Hashmaps
 - Locality-Sensitive Hashing (LSH)



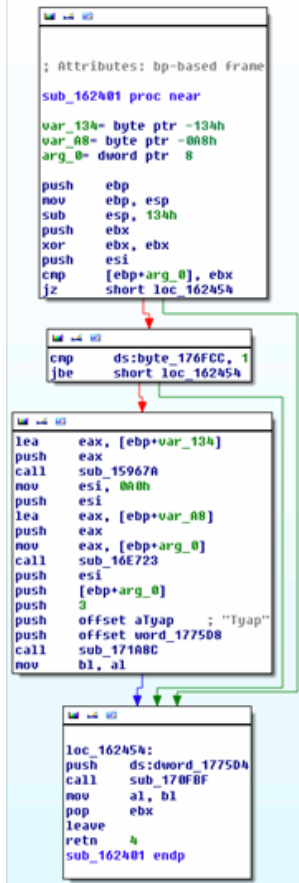
Project Co-Authors: Paul Hordienko, Steffen Enders, Manuel Blatt, Daniel Enders

[1] I. U. Haq and J. Caballero, "A survey of binary code similarity," In: Arxiv.org Computers & Security, 2019.

[2] S. Alrabaee, P. Shirani, L. Wang, and M. Debbabi, "FOSSIL: A Resilient and Efficient System for Identifying FOSS Functions in Malware Binaries," In: ACM Trans. Priv. Secur., vol. 21, 2018.

MCRIT

PIC Hashing



- Quasi-Identical: Position Independent Code (PIC) Hashing
 - On **function level** (original method as introduced by Cohen and Havrilla [1])
 - On basic block level (more granularity, almost the same speed)

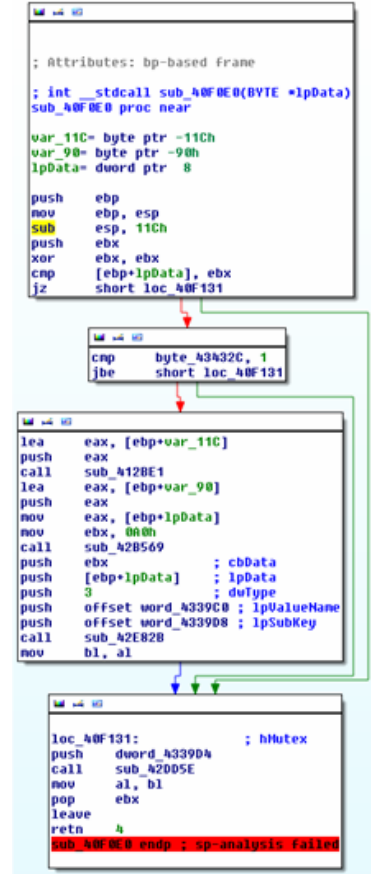
pichash(„55 8BEC 81EC**34**010000 ...“)

-> **8806641384121875405**

pichash(„55 8BEC 81EC**1C**010000 ...“)

-> **10270976525648996728**

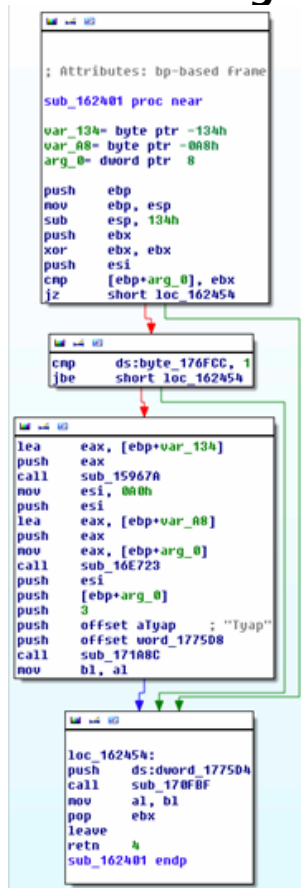
- Good for recognizing statically linked code (often binary identical)



[1] C. Cohen and J. Havrilla, "Function Hashing for Malicious Code Analysis", tech. rep., SEI, CMU, 2009.

MCRIT

PIC Hashing



■ Quasi-Identical: Position Independent Code (PIC) Hashing

- On function level (original method as introduced by Cohen and Havrilla [1])
- On **basic block level, size 4+** (more granularity, almost the same speed)

pi chash(„55 8BEC 81EC34010000 ...“)

-> 620962970

pi chash(„55 8BEC 81EC1C010000 ...“)

-> 2827249019

pi chash(„8D85CCFEFFFF 50 E8?? ...“)

-> 847901973

pi chash(„8D85E4FEFFFF 50 E8?? ...“)

-> 640583737

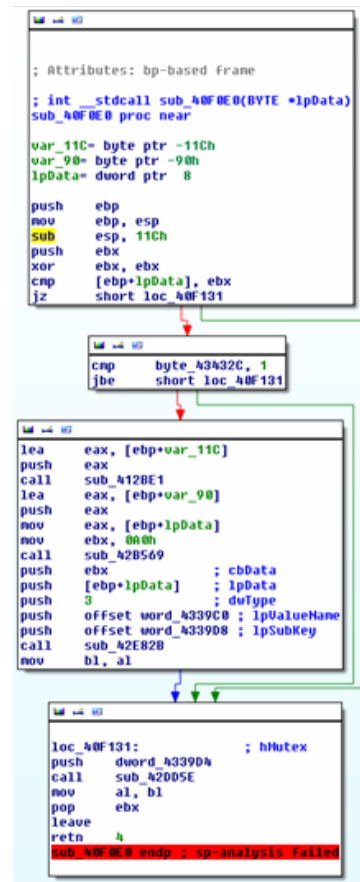
■ Good for recognizing partial code reuse (candidates)

pi chash(„FF35???????? E8???? ...“)

-> 463987246

pi chash(„FF35???????? E8???? ...“)

-> 463987246



[1] C. Cohen and J. Havrilla, "Function Hashing for Malicious Code Analysis", tech. rep., SEI, CMU, 2009.

MCRIT

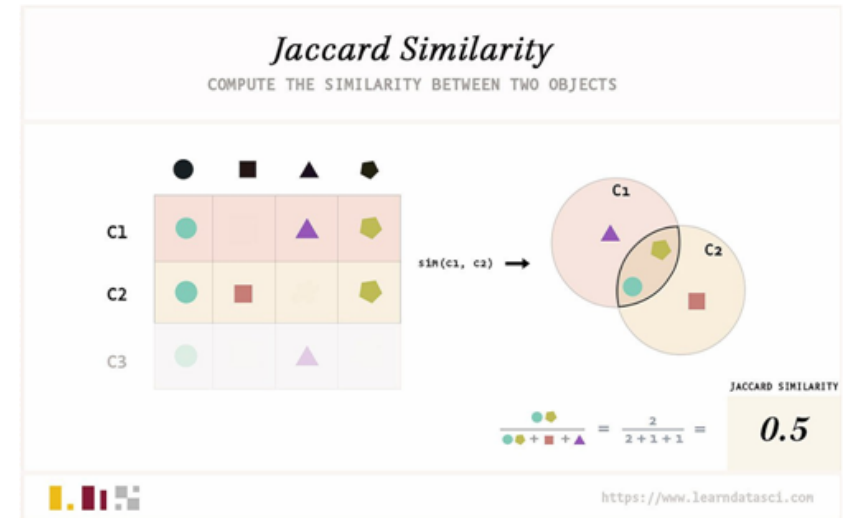
MinHash 101

■ MinHashing

- „Min-wise independent permutations“ - Locality Sensitive Hashing (LSH) scheme [1]
- Fast **estimation** of **set similarity** -> approximation of **Jaccard** similarity coefficient
- Scalability: $O(\log n)$ for single lookups

■ Use cases:

- text documents / websites (duplicates, plagiarism)
- genome sequencing
- code similarity! [2]



[1] "Min-wise independent permutations". Broder et al., In: Proceedings of the 30th ACM Symposium on Theory of Computing (STOC '98), New York, NY, USA.

[2] "Binary Function Clustering using Semantic Hashes". Jin et al., Carnegie Mellon University, 2012.

13 [3] <https://www.learn datasci.com/glossary/jaccard-similarity/>

MinHash Composition

Token-based features:

- Instruction 3-grams
- Abstract semantically
- Convert to 3-perms (sorted)

Metrics-based features:

- Numerically describe structure of a function
- Normalize & Quantize for fuzzy matching

MinHash matching:

- Count same values in same position

Example Function:

```

push esi
mov esi, 0x023C598C
push esi
call dword ptr [0x02391294]
push 0x14
pop eax
push dword ptr [0x023C5978]
mov word ptr [0x023C5980], ax
mov eax, dword ptr [esp+0x10]
mov dword ptr [0x023C5988], eax
call dword ptr [0x023B9085]
xor eax, eax
cmp dword ptr [esp+0x8], eax
jz 0x023B4830

```

```

push dword ptr [esp+0x8]
push eax
or eax, 0xFFFFFFFF
call 0x023B92C5

```

```

push esi
mov dword ptr [0x023C5978], eax
call dword ptr [0x0239129C]
pop esi
ret 0x8

```

Metrics-based Features

Feature	Value	LogBucket Triplets
max_block_size	14	12 14 16
num_calls	4	3 4 6
num_ins_C	7	4 6 8
num_ins_S	9	6 8 10
num_ins_A_rel	8	6 8 10
num_ins_M_rel	17	14 16 20
stack_size	0	-1 0 1

Token-based Features:

Instruction 3-grams (first 5)

```

push esi - mov esi, 0x023C598C - push esi
mov esi, 0x023C598C - push esi - call dword ptr [0x02391294]
push esi - call dword ptr [0x02391294] - push 0x14
call dword ptr [0x02391294] - push 0x14 - pop eax
push 0x14 - pop eax - push dword ptr [0x023C5978]

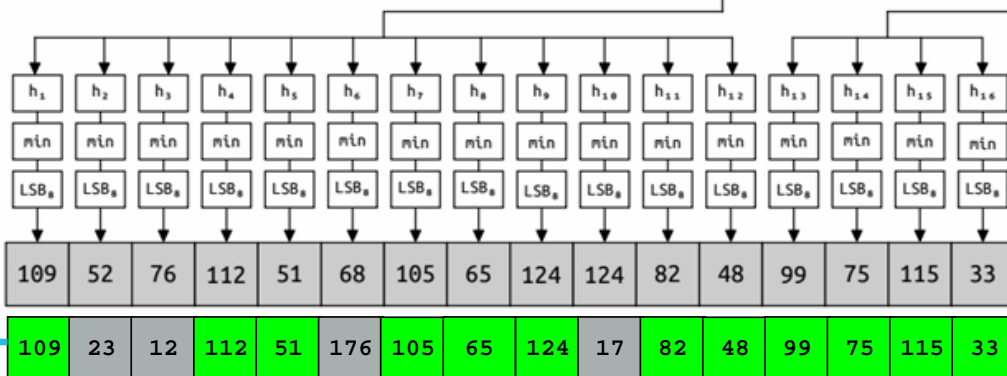
```

converted to escaped 3-perms:

```

M REG, CONST - S REG - S REG
C PTR - M REG, CONST - S REG
C PTR - S CONST - S REG
C PTR - S CONST - S REG
S CONST - S PTR - S REG

```



12/16 = 75%

MCRIT

MinHash Composition

Token-based features:

- Instruction 3-grams
- Abstract semantically
- Convert to 3-perms (sorted)

Metrics-based features:

- Numerically describe structure of a function
- Normalize & Quantize for fuzzy matching

Candidate Identification:

- Additionally index subsequences from signatures into buckets

Example Function:

```
push esi
mov esi, 0x023C598C
push esi
call dword ptr [0x02391294]
push 0x14
pop eax
push dword ptr [0x023C5978]
mov word ptr [0x023C5980], ax
mov eax, dword ptr [esp+0x10]
mov dword ptr [0x023C5988], eax
call dword ptr [0x023B9085]
xor eax, eax
cmp dword ptr [esp+0x8], eax
jz 0x023B4830
```

```
push dword ptr [esp+0x8]
push eax
or eax, 0xFFFFFFFF
call 0x023B92C5
```

```
push esi
mov dword ptr [0x023C5978], eax
call dword ptr [0x0239129C]
pop esi
ret 0x8
```

Metrics-based Features

Feature	Value	LogBucket Triplets
max_block_size	14	12 14 16
num_calls	4	3 4 6
num_ins_C	7	4 6 8
num_ins_S	9	6 8 10
num_ins_A_rel	8	6 8 10
num_ins_M_rel	17	14 16 20
stack_size	0	-1 0 1

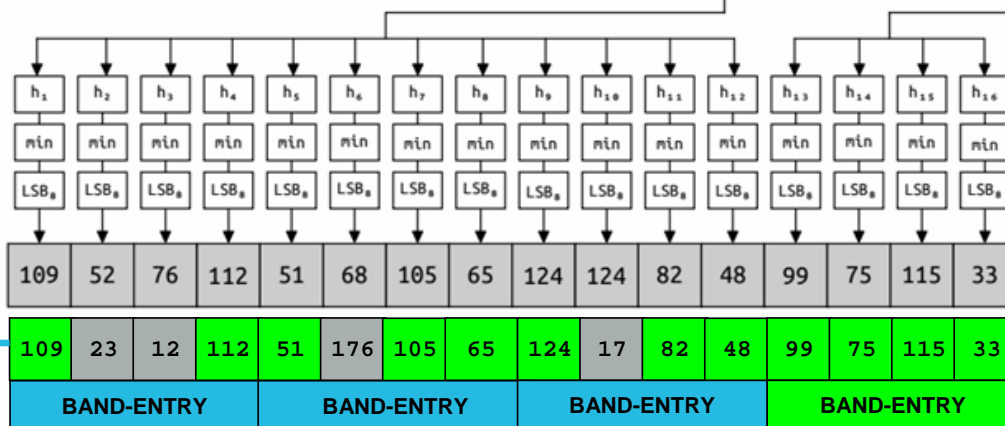
Token-based Features:

Instruction 3-grams (first 5)

```
push esi - mov esi, 0x023C598C - push esi
mov esi, 0x023C598C - push esi - call dword ptr [0x02391294]
push esi - call dword ptr [0x02391294] - push 0x14
call dword ptr [0x02391294] - push 0x14 - pop eax
push 0x14 - pop eax - push dword ptr [0x023C5978]
```

converted to escaped 3-perms:

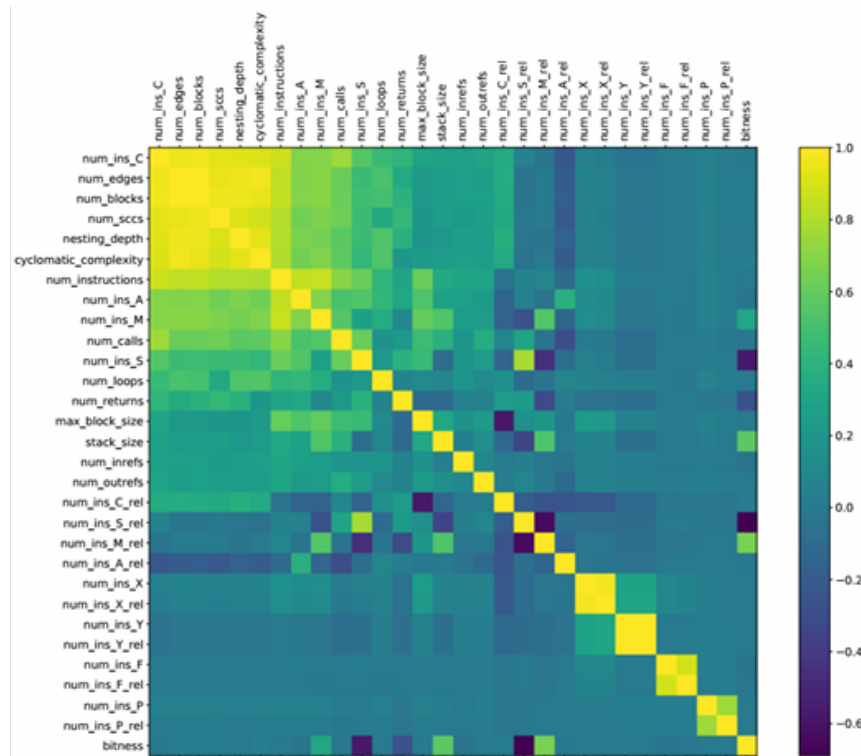
```
M REG, CONST - S REG - S REG
C PTR - M REG, CONST - S REG
C PTR - S CONST - S REG
C PTR - S CONST - S REG
S CONST - S PTR - S REG
```



12/16 = 75%

MinHash Feature Engineering

feature-name	space	min	25%	50%	75%	max	avg	sd	ρ_{32bit}	ρ_{64bit}	ρ_{same}	ρ_{diff}
cyclomatic_complexity	232	0	2	4	9	561	8.223	13.767	0.961	0.953	0.958	0.923
max_block_size	554	2	9	12	17	6,669	22.173	100.989	0.946	0.966	0.963	0.805
nesting_depth	65	0	2	4	7	101	4.871	4.777	0.973	0.966	0.970	0.960
num_blocks	339	1	4	9	18	817	15.417	22.605	0.974	0.939	0.954	0.933
num_calls	148	0	1	3	7	380	6.308	10.260	0.980	0.980	0.980	0.952
num_edges	480	0	4	11	25	1,376	21.640	36.096	0.969	0.949	0.955	0.929
num_inrefs	114	0	0	0	1	224	1.018	3.907	0.396	0.861	0.605	0.303
num_ins_A	516	0	3	6	12	3,929	16.244	83.531	0.985	0.963	0.981	0.688
num_ins_C	495	0	7	15	32	1,343	26.784	38.942	0.980	0.971	0.976	0.964
num_ins_F	14	0	0	0	0	30	0.003	0.161	0.972	-	0.972	-0.000
num_ins_M	691	0	6	14	33	2,811	29.725	60.236	0.980	0.970	0.981	0.384
num_ins_P	25	0	0	0	0	40	0.012	0.373	0.606	0.693	0.655	0.471
num_ins_S	301	0	2	6	15	932	13.854	24.591	0.947	0.946	0.965	0.006
num_ins_X	261	0	0	0	0	5,897	2.569	66.481	0.893	0.957	0.934	0.637
num_ins_Y	39	0	0	0	0	280	0.161	4.015	0.967	0.989	0.982	0.929
num_ins_A_rel	83	0	9	13	17	86	14.465	9.469	0.914	0.944	0.954	0.106
num_ins_C_rel	97	0	26	33	40	100	33.017	12.931	0.909	0.964	0.952	0.424
num_ins_F_rel	13	0	0	0	0	27	0.004	0.236	0.972	-	0.972	-0.000
num_ins_M_rel	100	0	18	31	46	99	32.194	16.829	0.935	0.939	0.975	-0.648
num_ins_P_rel	13	0	0	0	0	13	0.007	0.207	0.629	0.716	0.677	0.418
num_ins_S_rel	80	0	3	11	30	86	17.022	16.306	0.907	0.963	0.976	-0.339
num_ins_X_rel	97	0	0	0	0	99	1.068	7.057	0.887	0.957	0.933	0.636
num_ins_Y_rel	45	0	0	0	0	84	0.141	2.501	0.980	0.990	0.987	0.917
num_instructions	1,293	10	23	48	100	6,743	89.718	182.991	0.950	0.961	0.973	0.771
num_loops	36	0	0	0	1	100	0.380	1.139	0.957	0.962	0.961	0.894
num_outrefs	88	0	0	0	0	128	0.692	2.843	0.886	0.886	0.886	0.810
num_returns	56	0	1	2	3	232	2.197	3.249	0.973	0.930	0.956	0.682
num_scgs	262	1	4	8	15	555	12.672	16.598	0.973	0.941	0.954	0.931
stack_size	389	0	0	8	48	4,092	44.652	143.470	0.952	0.912	0.971	-0.292



Full details in my PhD thesis

Querying the System

■ Query with

- Basic Block
- Function
- Sample



	offset	num_bytes			Score
4459063 ▼	0x406770	150	win.hardrain	909	889162
4459063 ▼	0x406770	150	win.hardrain	1159	1135454
4459063 ▼	0x406770	150	win.op_blockbuster	1859	1816610
4459063 ▼	0x406770	150	win.op_blockbuster	5954	6681224
4459063 ▼	0x406770	150	win.romeos	4498	4990660
4459063 ▼	0x406770	150	win.romeos	4730	5238246
4459063 ▼	0x406770	150	win.romeos	5710	6362927
4459063 ▼	0x406770	150	win.badcall	1288	1257856

For every function, we know **how many** families we match.

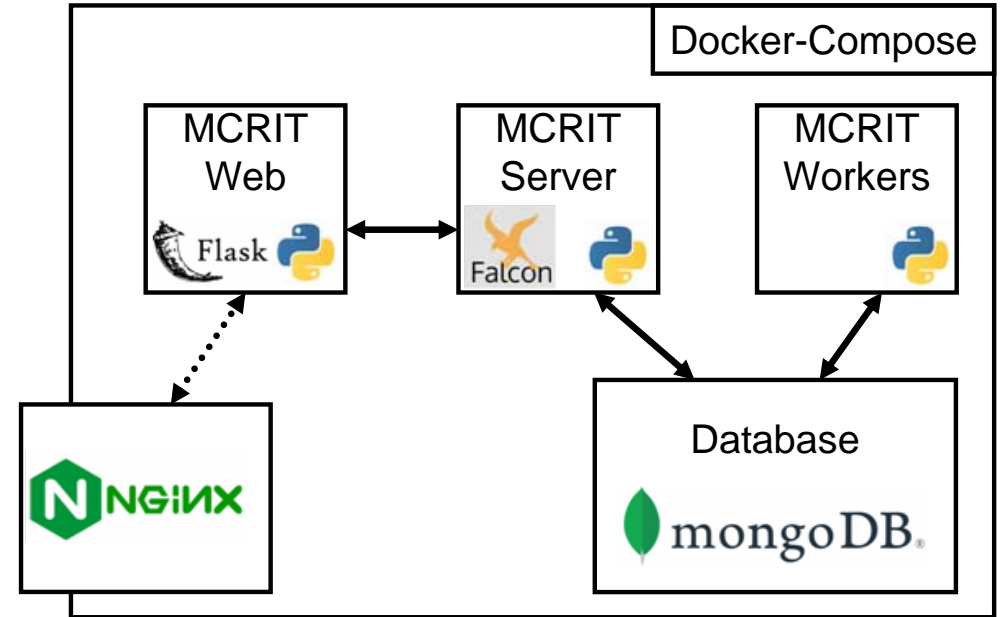
1. We can use this to weigh by occurrence **frequency** when aggregating to sample matches.
2. We can identify **unique** matches into just one family and use this as further indication for identity.

	Version		SHA256	Filename	Bitness	FNs	Min#	Pic#	Lib	Direct	Frequency		Uniq
win.romeos ▼	2014-07-07-alfa	4498 ▼	4ec8214b	eff542ac8e...0x00400000	32	311	247	247	71	98	47	59	21.49%
win.hardrain ▼		909 ▼	c3b1af35	2cc3b5f2df...0x00340000	32	289	138	103	66	50	13	13	0.00%
win.op_blockbuster ▼	2015-04-08	1859 ▼	2f629c3c	2f629c3c65...3_unpacked	32	341	159	122	70	56	12	12	0.00%
win.keymarble ▼	2017-04-12	4543 ▼	f19cd9ef	e23900b00f...0x00400000	32	448	168	137	70	53	9	8	0.00%

[1] eff542ac8e37db48821cb4e5a7d95c044ff27557763de3a891b40eb52cc55 @ 0x406770 | win.romeos |

MCRIT Setup

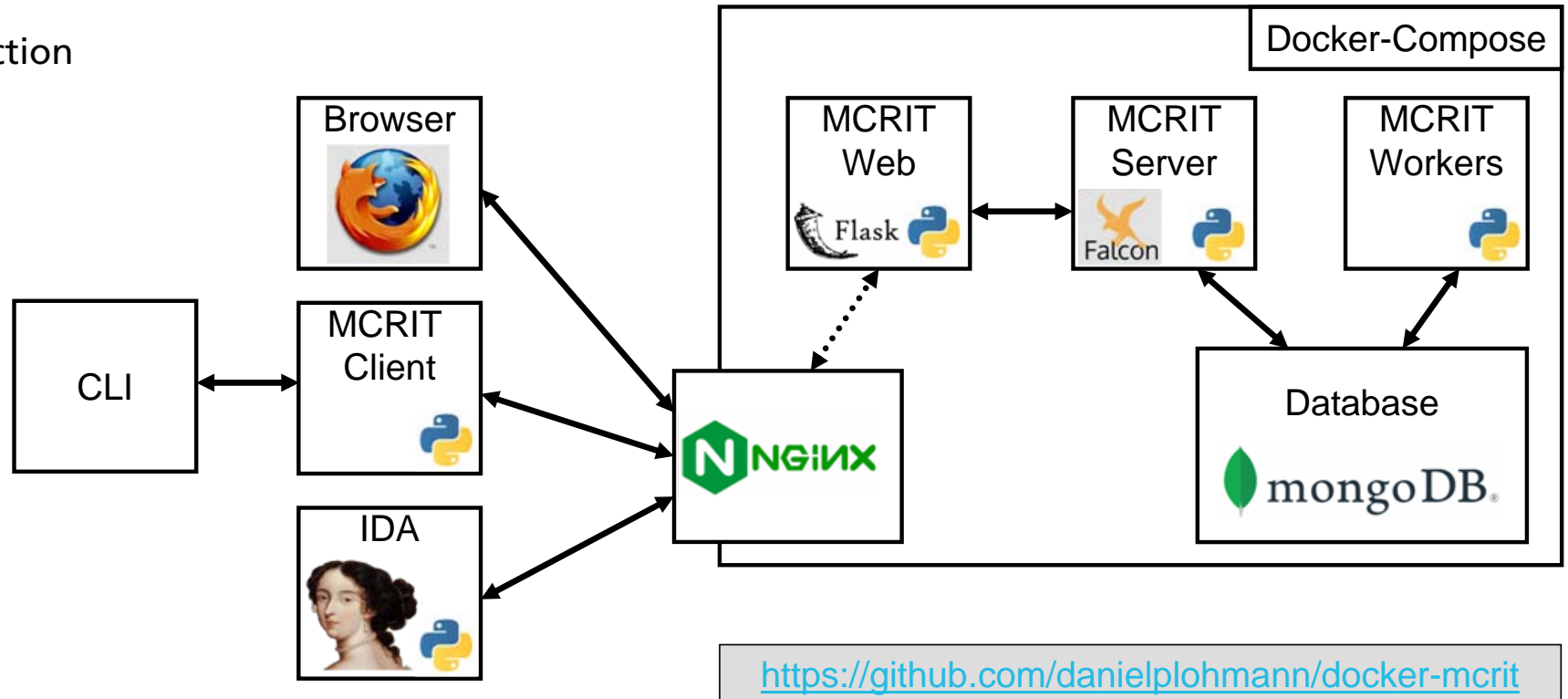
- Database
- MCRIT Server
 - Core of the system
 - Enable access to stored content (API)
 - Create matching jobs
- MCRIT Workers
 - Process jobs from the queue
- MCRIT Web
 - Expose service functionality in a user interface
 - User management
 - API forwarding to MCRIT server



<https://github.com/danielplohm/docker-mcrit>


MCRIT Setup


■ Interaction



MCRIT

Setup: WebUI


































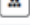

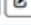
 **MCRIT**

 **Fraunhofer**
FKIE

Explore Analyze Data pnx

Overview of Families

Search

★	Family	Samples	Functions	Library	
0	Unnamed 🔗	38	55170	🔴	   
1	msvcrt 🔗	327	310944	🟢	   
2	win.wastedlocker 🔗	10	1620	🔴	   
3	win.isfb 🔗	118	74792	🔴	   
4	win.juicy_potato 🔗	1	1688	🔴	   
5	win.lyposit 🔗	2	1095	🔴	   
6	win.horus_eyes_rat 🔗	0	0	🔴	   
7	win.bumblebee 🔗	23	157017	🔴	   
8	win.cobra 🔗	53	14954	🔴	   

MCRIT

Setup: WebUI

The screenshot displays the MCRIT WebUI interface. At the top, there is a header with the MCRIT logo and the Fraunhofer FKIE logo. Below the header, there are navigation tabs: 'Explore', 'Analyze', 'Data', and 'pnx'. The main section is titled 'Compare Samples 1vsN'. A search bar contains the text 'wannacry'. Below the search bar, there is a table with the following columns: SHA256, Family, Version, Filename, Bitness, Functions, and Library. The table contains six rows of sample data. At the bottom of the interface, there is a 'Minhash Matching' section with a 'Force rematch' checkbox and a slider. A 'Compare' button is located at the bottom right.

SHA256	Family	Version	Filename	Bitness	Functions	Library
626	e458d473 win.wannacryptor	vt-2017-05-05	0345782378ee7a8b48c2...7366_dump_0x00400000	32	922	
1380	2be58051 win.wannacryptor	2017-02-09	3e6de9e2baacf9309496...eed9_dump_0x00400000	32	450	
4092	ca29de1d win.wannacryptor	2017-03-19	ca29de1dc8817868c93e...1f52ba469c8_unpacked	32	907	
4805	d181360a win.wannacryptor	vt-2017-05-12	b9c5d4339809e0ad9a00...1c25_dump_0x00400000	32	926	
4955	d36a4116 win.wannacryptor	vt-2017-05-12	ed01ebfbc9eb5bba545...41aa_dump_0x00400000	32	926	
5828	6611cc9e win.wannacryptor	2017-03-19	ca29de1dc8817868c93e...69c8_dump_0x00400000	32	590	

Best Family Matches

total: 5, showing: 1 - 5 (filtered: 966)

Filter results to (nonlib) direct score: regular (0-100) nonlib (0-100)

Filter results to (nonlib) frequency score: 5 nonlib (0-100)

☐ only show families with unique matches ☐ exclude own family

[filter](#) [clear](#)

★	Version	★	SHA256	Filename	Bitness	FNs	Min#	Pic#	Lib	Direct	Frequency	Uniq
win.wannacryptor ▼	2017-03-19	4092 ▼	ca29de1d	ca29de1dc8...8_unpacked	32	907	161	114	14	69 68	34 34	13.93%
win.kuaibu8 ▼	2017-01-30	3490 ▼	3d3c3bf0	5793b30729...0x00400000	32	1288	76	62	8 35	34	9 9	0.00%
win.sys10 ▼	2013-03-07	4357 ▼	6386ae55	afe3dd68bd...0x00400000	32	337	48	37	5 31	31	8 8	0.00%
win.joanap ▼		5695 ▼	31e27637	a1c483b0ee...0x002c0000	32	186	30	24	3 11	11	7 7	0.00%
win.badcall ▼		2923 ▼	d594b683	93e13ffd2a...0x006e0000	32	166	58	44	3 27	27	7 7	0.00%

« < 1 > »

[Compare](#)

Use Cases

Example Use Cases

- Malware family identification and library code differentiation
- Isolation of unique family code
- Lead generation for discovering potentially unknown links
- Label Transfer

Example Use Cases

Malware Family Identification and Library Code Differentiation

LockBit ransomware goes 'Green,' uses new Conti-based encryptor

By [Lawrence Abrams](#)

February 1, 2023 05:48 PM



The LockBit ransomware gang has again started using encryptors based on other operations, this time switching to one based on the leaked source code for the Conti ransomware.

Since its launch, the LockBit operation has gone through numerous iterations of its encryptor, starting with a custom one and moving to LockBit 3.0 (aka LockBit Black), which is derived from the BlackMatter gang's source code.

This week, cybersecurity collective VX-Underground first reported that the ransomware gang is now using a new encryptor named 'LockBit Green,' based on the leaked source code of the now-disbanded Conti gang.

num_samples	6243
num_families	1440
num_functions	6995154
num_pichashes	1464810

Query Sample

Drop file or click here to import

45c317_lockbit_green.exe

☒ Unmapped
☐ Dumped

Minhash Matching: Standard

Submit

Disassembly + Matching: 35sec

[1] <https://www.bleepingcomputer.com/news/security/lockbit-ransomware-goes-green-uses-new-conti-based-encryptor/>

[2] lockbit green: 45c317200e27e5c5692c59d06768ca2e7eeb446d6d495084f414d0f261f75315

Best Family Matches

total: 1070, showing: 1 - 10

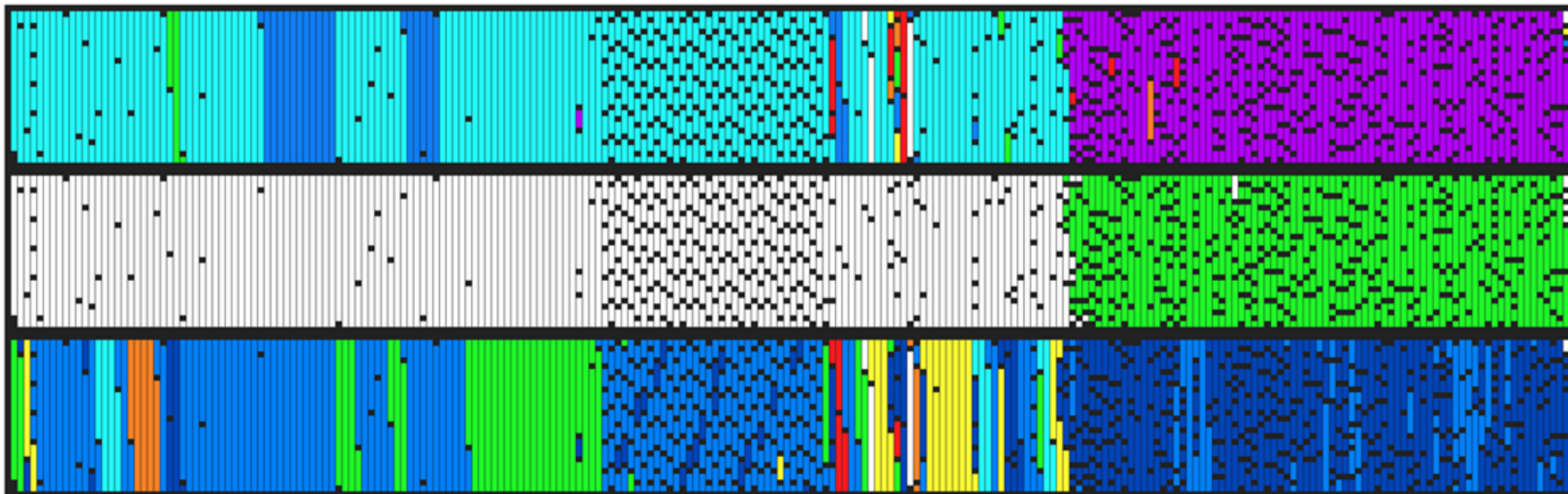
🌟	Version	🌟	SHA256	Filename	Bitness	FNs	Min#	Pic#	Lib	Direct	Frequency		
win.meow ▼		2172 ▼	222e2b91	222e2b91f5...3_unpacked	32	702	461	126	220	66	71	41	53
win.conti ▼	2021-02-04	5041 ▼	a5751a46	a5751a4676...6_unpacked	32	736	516	183	256	59	58	28	36
win.scarecrow ▼		6199 ▼	bcf49782	bcf49782d7...a_unpacked	32	653	582	271	334	61	51	27	32
win.lockergoga ▼	2019-03-18	4517 ▼	edae201c	c97d9bbc80...0x00400000	32	7847	369	311	352	26	1	3	0
win.void ▼		1460 ▼	2fd1863e	2fd1863eb3...c_unpacked	32	7123	366	312	351	26	1	3	0
win.bandook ▼		4792 ▼	fabce973	fabce973a9...7_unpacked	32	3229	363	306	349	25	1	3	0

Best Family Matches

total: 1070, showing: 1 - 10

🌟	Version	🌟	SHA256	Filename	Bitness	FNs	Min#	Pic#	Lib	Direct	Frequency		
win.meow ▼		2172 ▼	222e2b91	222e2b91f5...3_unpacked	32	702	461	126	220	66	71	41	53
win.conti ▼	2021-02-04	5041 ▼	a5751a46	a5751a4676...6_unpacked	32	736	516	183	256	59	58	28	36
win.scarecrow ▼		6199 ▼	bcf49782	bcf49782d7...a_unpacked	32	653	582	271	334	61	51	27	32

Showing: foreign family match frequency, library matches, best foreign family match scores.



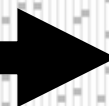
Best Family Matches

total: 1070, showing: 1 - 10

★	Version	★	SHA256	Filename	Bitness	FNs	Min#	Pic#	Lib	Direct	Frequency	
win.meow ▼		2172 ▼	222e2b91	222e2b91f5...3_unpacked	32	702	461	126	220	66 71	41	53
win.conti ▼	2021-02-04	5041 ▼	a5751a46	a5751a4676...6_unpacked	32	736	516	183	256	59 58	28	36
win.scarecrow ▼		6199 ▼	bcf49782	bcf49782d7...a_unpacked	32	653	582	271	334	61 51	27	32

Showing: foreign family match frequency, library matches, best foreign family match scores.

Code from standard libraries (MSVCRT, ...)

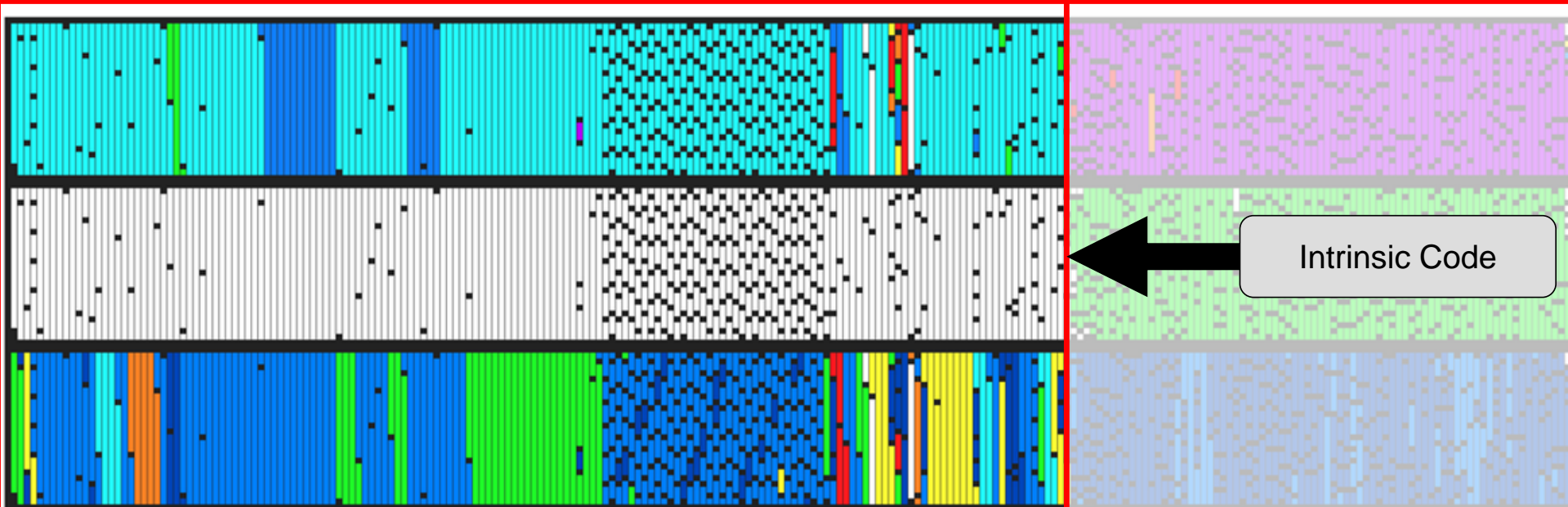


Best Family Matches

total: 1070, showing: 1 - 10

★	Version	★	SHA256	Filename	Bitness	FNs	Min#	Pic#	Lib	Direct	Frequency	
win.meow ▼		2172 ▼	222e2b91	222e2b91f5...3_unpacked	32	702	461	126	220	66 71	41	53
win.conti ▼	2021-02-04	5041 ▼	a5751a46	a5751a4676...6_unpacked	32	736	516	183	256	59 58	28	36
win.scarecrow ▼		6199 ▼	bcf49782	bcf49782d7...a_unpacked	32	653	582	271	334	61 51	27	32

Showing: foreign family match frequency, library matches, best foreign family match scores.



Best Family Match

total: 1070, showing: 1 - 10

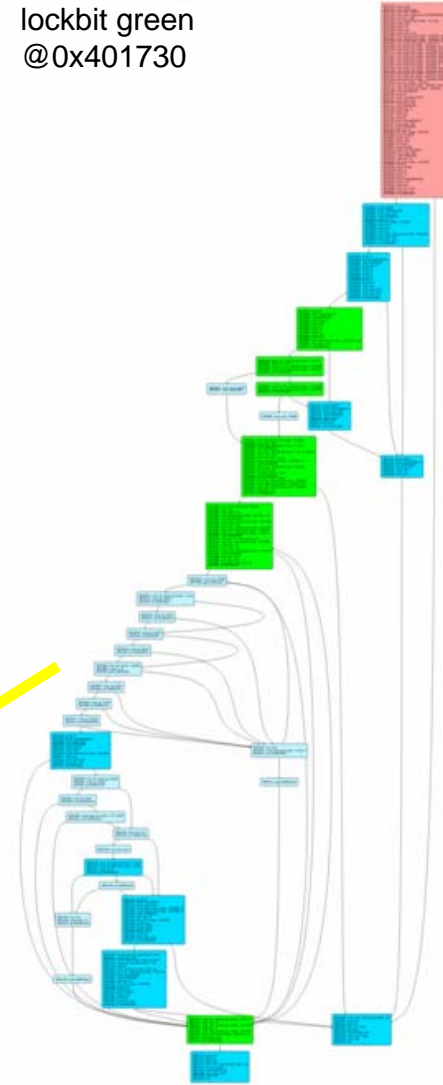
★	Version
win.meow ▼	
win.conti ▼	2021-02-
win.scarecrow ▼	

lockbit green
@0x401730

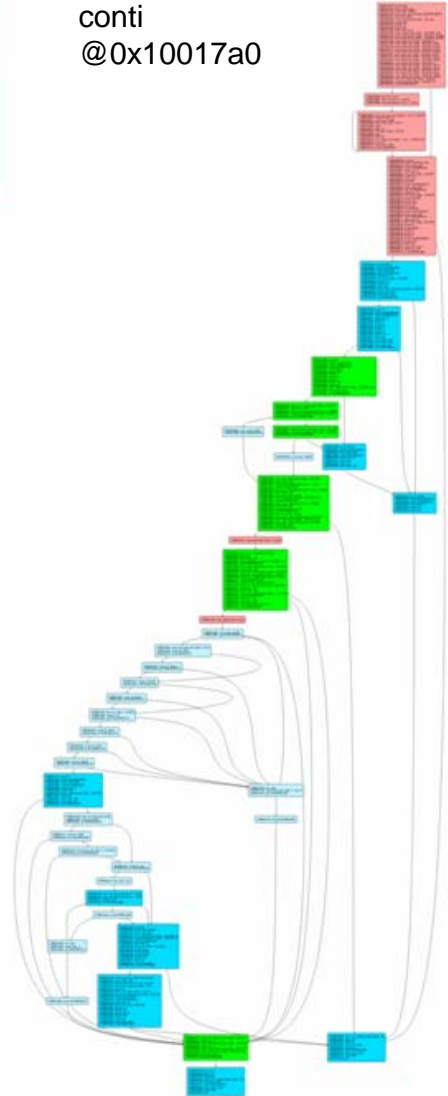
conti
@0x10017a0

Function CFGs

Show Cycles Show Loops ☒ Show Loop Boundaries ☐ Enable Tooltip

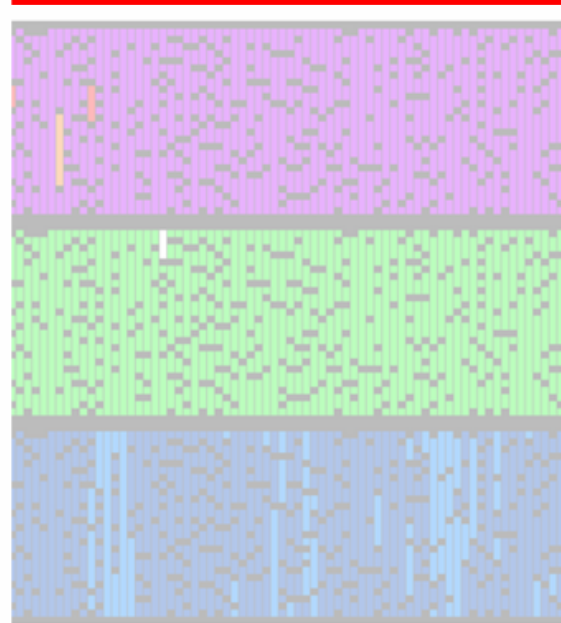
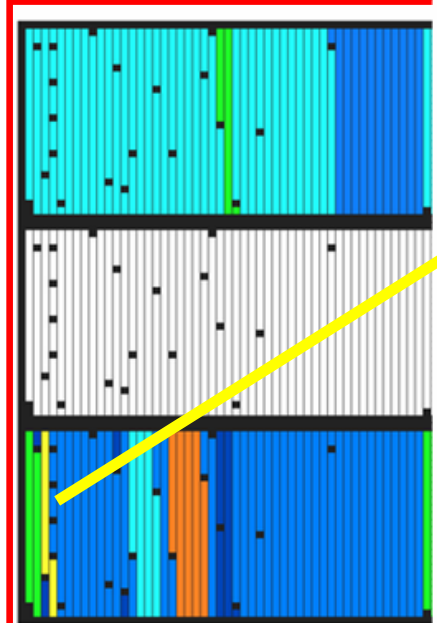


Go to top



Min#	Pic#	Lib	Direct	Frequency
461	126	220	66 71	41 53
516	183	256	59 58	28 36
582	271	334	61 51	27 32

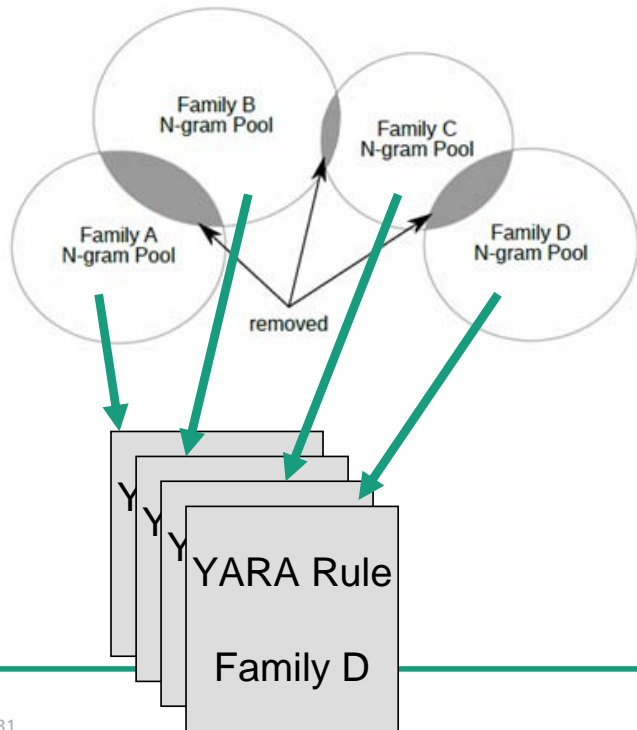
Showing: foreign family match frequency



Example Use Cases

Isolation of Unique Family Code

- Essentially like YARA-Signator, but with basic blocks



Unique Block Isolation Report

Job ID	639358a4e0ff5413a77221e4
Family	win.remcos
Samples	19
Unique Blocks	10028
Has a YARA rule?	True, covers: 19 samples
YARA rule covers all?	True

Statistics Unique Blocks YARA Rule

Block Statistics across Samples

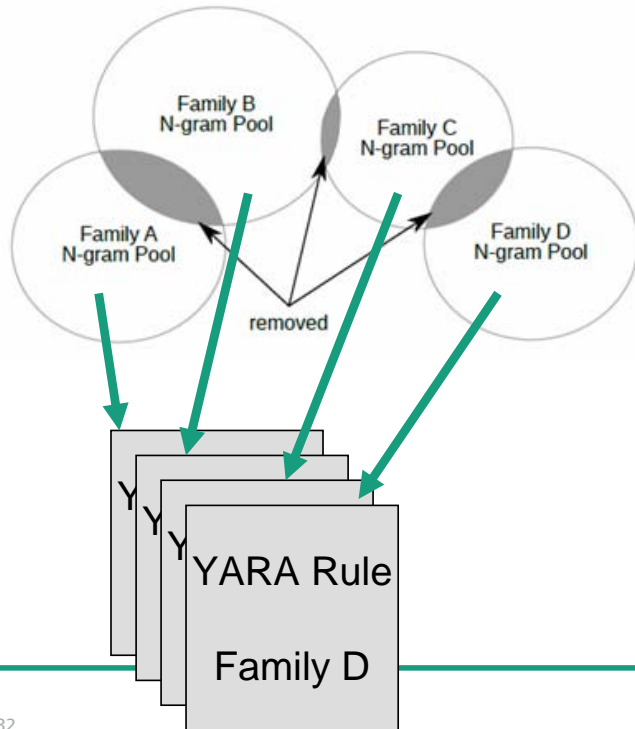
Characteristic blocks are basic blocks only found in this collection of samples (versus rest of the whole data set), unique blocks are only found in the specific sample.

Sample ID	Total Blocks	Characteristic Blocks	Unique Blocks
1111	1104	965 (87.41%)	0 (0.00%)
1342	1178	1056 (89.64%)	567 (48.13%)
1343	1106	967 (87.43%)	1 (0.09%)
1519	746	548 (73.46%)	68 (9.12%)
1968	8020	3423 (42.68%)	188 (2.34%)
2211	867	743 (85.70%)	266 (30.68%)
3729	1190	1041 (87.48%)	32 (2.69%)
4501	8307	3937 (47.39%)	0 (0.00%)
4561	8315	3938 (47.36%)	0 (0.00%)
4575	1299	1137 (87.53%)	74 (5.70%)
4683	876	745 (85.05%)	49 (5.59%)

Example Use Cases

Isolation of Unique Family Code

- Essentially like YARA-Signator, but with basic blocks



Explore Unique Blocks

Filter blocks to

filter

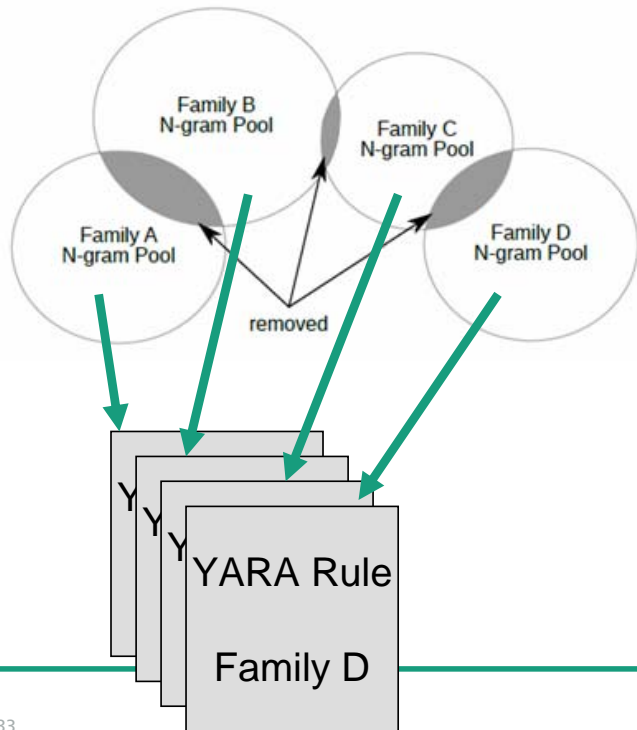
total: 10028, showing: 1 - 100

Score	PicBlockHash	Samples	Instructions	Function ID	Block
97.14	0xcadb40567ec2656	19 / 19	6	470589	<pre>/* picblockhash: 0xcadb40567ec2656 * 6a09 push 9 * ff35fc804100 push dword ptr [0x418dfc] * ff3554244100 call dword ptr [0x412454] * ff35fc804100 push dword ptr [0x418dfc] * ff3508244100 call dword ptr [0x412488] * e038 jmp 0x41157f */ [6a09 ff3577777777 ff3577777777 ff3577777777 ff3577777777 e077]</pre>
83.16	0x1bb7e1ed48bd751	15 / 19	7	470361	<pre>/* picblockhash: 0x1bb7e1ed48bd751 * 53 push ebx * 53 push ebx * 56 push esi * 68ed524000 push 0x4052ed * 53 push ebx * 53 push ebx * ffd7 call edi */ [53 53 56 6877777777 53 53 ffd7]</pre>
80.30	0x38bb1b656ac38756	15 / 19	6	470428	<pre>/* picblockhash: 0x38bb1b656ac38756 * 6890374100 push 0x413790 * 6874374100 push 0x413774 * ffd7 call edi * 50 push eax * ffd6 call esi * a3a48a4100 mov dword ptr [0x418aa8], eax */ [6877777777 6877777777 ffd7 50 ffd6 a377777777]</pre>

Example Use Cases

Isolation of Unique Family Code

- Essentially like YARA-Signator, but with basic blocks



Proposed YARA rule

Copy rule to clipboard!

```
rule mcrit_639358a4e0ff5413a77221e4 {
  meta:
    author = "MCRIIT YARA Generator"
    description = "Code-based YARA rule composed from potentially unique basic blocks for the selected set of samples/family."
    date = "2023-03-29"

  strings:
    // Rule generation selected 20 picblocks, covering 19/19 input sample(s).
    /* picblockhash: 0x1b07e1e048b1d751 - coverage: 15/19 samples.
    * 53      | push ebx
    * 53      | push ebx
    * 56      | push esi
    * 68ed524000 | push 0x4052ed
    * 53      | push ebx
    * 53      | push ebx
    * ffd7    | call edi
    */
    $picblockhash_0x1b07e1e048b1d751 = { 53 53 56 68???????? 53 53 ffd7 }

    /* picblockhash: 0x2e684020db1f147c - coverage: 14/19 samples.
    * ff7508  | push dword ptr [ebp + 8]
    * 6a00    | push 0
    * 6800040000 | push 0x400
    * ff15cc204100 | call dword ptr [0x4120cc]
    * 804008    | test ecx, [ebp + 8]
    * 51        | push ecx
    * 50        | push eax
    * ff15e4804100 | call dword ptr [0x4180e4]
    * 83700800    | cmp dword ptr [ebp + 8], 0
    * 7504       | jne 0x418ae5
    */
    $picblockhash_0x2e684020db1f147c = { ff7508 6a00 6800040000 ff15???????? 804008 51 50 ff15???????? 83700800 75?? }

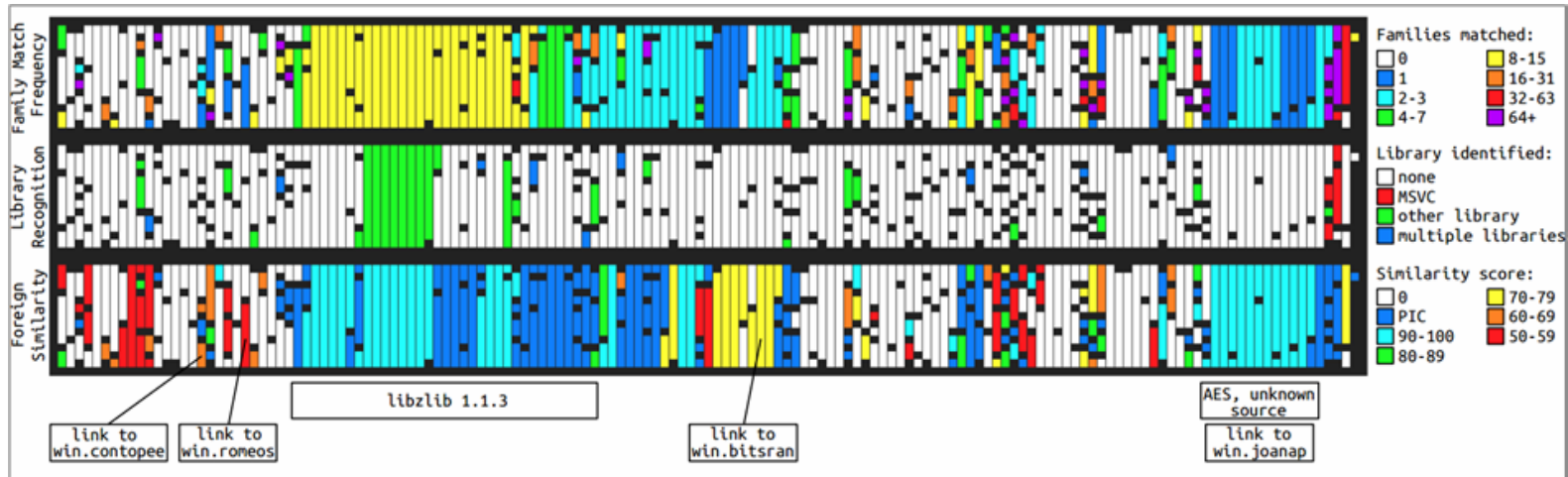
    /* picblockhash: 0x380b1d656ac38756 - coverage: 15/19 samples.
    * 6890374100 | push 0x413790
    * 6874374100 | push 0x413774
    * ffd7      | call edi
    * 50        | push eax
    * ffd6      | call esi
    * a3a88a4100 | mov dword ptr [0x418aa4], eax
    */
    $picblockhash_0x380b1d656ac38756 = { 68???????? 68???????? ffd7 50 ffd6 a3???????? }

    /* picblockhash: 0x3fc01ffe2622434d - coverage: 14/19 samples.
    * ff7508  | push dword ptr [ebp + 8]
    * f796    | push dword ptr [esi]
    * ff15ac234100 | call dword ptr [0x4123ac]
    * 59      | pop ecx
    * 85c0    | test eax, eax
    * 59      | pop ecx
    * 741f    | je 0x40f2dc
    */
  }
```

Example Use Cases

Lead Generation for Discovering Potentially Unknown Links

- May 15th 2017 - Tweet by Neel Mehta (Google) with hashes + offsets
 - Earlier version of WannaCry sharing „rare“ code with Contopee
- Identification of similar functions with appearance across few families
 - Potential reuse of non-public code as an indicator for relationship



[1] <https://twitter.com/neelmehta/status/864164081116225536>

[2] „Classification, Characterization, and Contextualization of Windows Malware using Static Behavior and Similarity Analysis“, D. Plohmann, 2022.

Example Use Cases

Lead Generation for Discovering Potentially Unknown Links

Function Matches

selection: 49, showing: 1 - 49 (filtered: 128)

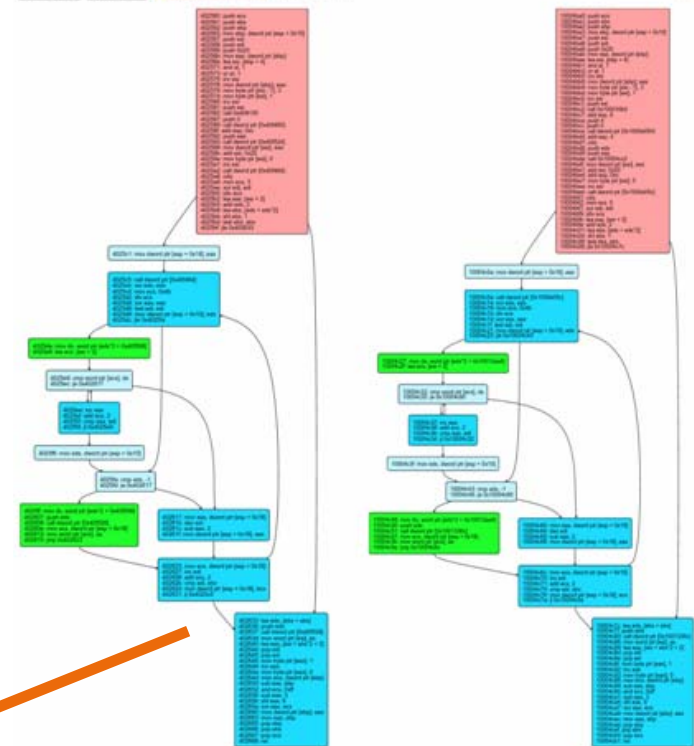
Filter results to max score (0-100)

☒ exclude functions with library hits ☐ exclude PIC hits

Function ID	offset	num_bytes	Matched Families	Matched Samples	Matched Functions	Best Score	Min	Pic	Lib	Uniq
1344938	0x401000	58	1 查	5	5	100	5	5	0	☑
1344939	0x401040	173	1 查	3	3	100	3	3	0	☑
1344941	0x401120	252	1 查	3	3	100	3	3	0	☑
1344942	0x401220	69	1 查	3	3	100	3	3	0	☑
1344948	0x401360	61	1 查	5	5	100	5	3	0	☑
1344953	0x401790	278	1 查	5	5	90	5	0	0	☑
1344955	0x4018d0	183	1 查	5	5	68	5	0	0	☑
1344961	0x401dd0	42	1 查	5	5	100	5	5	0	☑
1344965	0x401e40	145	1 查	5	5	100	5	5	0	☑
1344966	0x401ee0	190	1 查	5	5	100	5	5	0	☑
1344967	0x401fa0	181	1 查	5	5	84	5	0	0	☑
1344968	0x402060	344	1 查	5	5	84	5	0	0	☑
1344969	0x4021c0	77	1 查	5	5	78	5	0	0	☑
1344970	0x402210	451	1 查	5	5	95	5	0	0	☑
1344971	0x4023e0	214	1 查	5	5	82	5	0	0	☑
1344975	0x402500	55	3 查 win.contopee win.op_blockbuster win.roncoos	7	7	100	7	7	0	☑
1344977	0x402560	265	1 查	1	1	67	1	0	0	☑

Function CFGs

☒ Show Loop Boundaries ☐ Enable Tooltip



Example Use Cases

Label Transfer

MCrit4IDA v0.1

Activity Info: 2023-03-29T13-21-33Z - Success! Received remote Sample Entry.
Remote server: http://127.0.0.1:8000/ - 0.19.0 - No statistics.
Remote sample: 2 (win.wannacry -)

SHA256: 3e6de9e2baac930949647c399818e7a2caea2626df6a468407854aaa515eed9
Architecture: intel
Bitness: 32 bit
ImageBase: 0x401000
Functions: 428 (leaf: 222, recursive: 2)
Instructions: 17430
Code Size: 52406 bytes
Family: win.wannacry
Version:
Library: NO

Function Matches Sample Match Summary Function Match Summary

Matches for Function: 0x402560 - 2 families, 2 samples, 2 functions.
☐ Filter out Library Matches
☒ Active Live PicHash Queries

Function Matches

ID	SHA256	Sample	Family	Version	Pic#	Min#	Lib
1 1158	3e6de9e2	2	win.wannacry		YES	100	NO
2 85	2c4e5ba7	0	win.contopee		NO	67	NO

Names from Matched Functions

ID	Score	user	Function Label
1 1158	100	anonymous	init_cnc_packet

```
push    eax, [ebp+0]
lea     esi, [ebp+4]
and     al, 1
or      al, 1
inc     esi
mov     [ebp+0], eax
mov     byte ptr [esi-1], 3
mov     byte ptr [esi], 1
inc     esi
push    call
push    sub_408130
; Time
xor     move
smth_time
push    sub_4023E0
push    sub_4024C0
mov     mov
add     esi, 20h
mov     byte ptr [esi], 0
inc     esi
call    ds:rand
cdq
mov     ecx, 0
xor     edi, edi
idiv    ecx
lea     eax, [esi+0]
add     edx, 2
lea     ebx, [edx+edx*2]
shl     ebx, 1
test    ebx, ebx
jle     short loc_402633

mov     [esp+14h+arg_0], eax

loc_4025C3:
call    ds:rand
xor     edx, edx
mov     ecx, 4Bh
div     ecx
xor     eax, eax
test    edi, edi
mov     [esp+10h+var_4], edx
jle     short loc_4025FA
```

100.00% (~117,238) (928,740) 0000259E 0040259E: init_cnc_packet+3E (Synchronized with Hex View-1)

Output

```
Python>
Python>
init_cnc_packet
Python
```

Summary

Limitations

Minhash-based Code Relationship & Investigation Toolkit (MCRIT)

- Version released today is a first version
 - Fully functional, but needs some usability improvements
- Data exchange
 - Basic import / export, looking to improve reference data distribution
- Architecture support
 - x86/x64 only, not optimized for cross-bitness
- Matching / Search
 - Currently only PicHash and MinHash, may add further options (WinAPI, strings, PE/ELF meta data, ...)

[1] <https://github.com/danielplohmann/docker-mcrit>

Summary

Minhash-based Code Relationship & Investigation Toolkit (MCRIT)

■ MCRIT

- A framework for quasi-identical and fuzzy 1:n code matching

■ Variety of Use Cases

- Code identification & library filtering, hunting, label transfer, ...

■ Full Open Source Release

- Convenient deployment via Docker [1]



[1] <https://github.com/danielplohmann/docker-mcrit>

Thank you for your attention!

Dr. Daniel Plohmann
daniel.plohmann@fkie.fraunhofer.de

 @push_pnx

