Code Cartographer's Diary





2018-10-02 | One Conference, Den Haag

Daniel Plohmann daniel.plohmann@fkie.fraunhofer.de

Steffen Enders steffen.enders@tu-dortmund.de

Martin Clauß martin.clauss@fkie.fraunhofer.de

Paul Hordiienko pavlo.hordiienko@fkie.fraunhofer.de

Elmar Padilla elmar.padilla@fkie.fraunhofer.de





The Agenda



Agenda

Overview: Malpedia

- Motivation & Overview
- Operationalizing Malpedia
- Windows API Usage Recovery & Analysis for Malware Characterization
 - Tools: ApiScout / ApiVectors
 - Evaluation Results
- Code-based Similarity Analysis
 - Tools: SMDA & MCRIT
 - Evaluation Outlook
- Summary

3



malpedia



Δ

Overview What is Malpedia?



[1] https://malpedia.caad.fkie.fraunhofer.de 5 [2] https://malpedia.io

© Cyber Analysis and Defense Department, Fraunhofer FKIE

Overview Statistics



- Launched @ Botconf 12/2017 [3]
 - ~800 users, 150 monthly active
- 2000+ user contributions since then
 - THX to "you know who you are"!! :)

- [1] https://malpedia.caad.fkie.fraunhofer.de
- 6 [2] https://malpedia.io
- [3] https://www.botconf.eu/wp-content/uploads/2017/12/2017-DanielPlohmann-Malpedia.pdf
- © Cyber Analysis and Defense Department, Fraunhofer FKIE



Overview What does it look like?

malpedia		Fraunhofer	
_malpedia	a	Fraunhofer	
malped	lia	Inventory Statistics Usage Users 10 pnx	
The followin.citadel (Back to overview) 1 page III Citadel 2 page III Citadel 3		Propose Change	Alternatively:
There is no description at this point.			Or or
6 A https://blog.malwarebytes.com/threat	t-analysis/2012/11/citadel-a-cyber-criminals-ultimate-weapon	1	git cione malpedia
7 A http://www.xylibox.com/2016/02/cita	del-0011-atmos.html		to retrieve all data for
⁶ http://blog.jpcert.or.jp/2016/02/bankii	ng-trojan27d6.html		offline use at once
6 https://www.arbornetworks.com/blog 10	/asert/the-citadel-and-gameover-campaigns-of-5cb682c10440	b2ebaf9f28c1fe438468/	
Download Yara-Signature			
Samples			
Version	SHA256	Status VT	
2012-02-20-1.3.0.0	d147e69bd47cd504ca555fcad15081ac5bd79b584	1794a93d5b250eb73d05b186 dumped 🛯 🛓	



[1] <u>https://malpedia.caad.fkie.fraunhofer.de</u> 7 [2] <u>https://malpedia.io</u>

Overview

Operationalizing Malpedia

Identification

- YARA
- Search / Comparison
- Label Provider (Clustering)
- Contextualization
 - Publication references for families, actors, ...
- **QA** / Regression Testing
 - Tools, Config extractors, etc







Malpedia Usage Example:

Windows API Usage Recovery & Analysis for Malware Characterization



9

Motivation

"(Windows) API interactions are an essential cornerstone for effective reverse engineering"



Current "State of the Art"



Overview

Tool: ApiScout [1]

- Originally presented at Botconf, December 2017
- Library for painless (Windows) API reconstruction in known environments
- Idea: API function offset bruteforcing based on databases
- Extension: ApiVectors
 - Compact representation (bit vector) indicating the presence of relevant WinAPI functions
 - Enables fast assessment of malware's potential capabilities
 - Allows similarity analysis based on WinAPI usage characteristics



[1] https://github.com/danielplohmann/apiscout

ApiScout Methodology



Accuracy Evaluation:

Sample: 15 Windows Bins 5,367 API Imports

FKIF

WinAPI Availability for Static Analysis / Methods of API Usage



- Across 702 families (90 ignored -> .net)
- PE Imports:
 - From PE Header Import Table only
- Dynamic + Cached:
 - LoadLibrary / GetProcAddress ApiHashing -> Custom IAT
- Obfuscation:
 - Custom Jump Table (Andromeda)
 Offset-based Hook Avoidance (Chthonic)
 On-Demand Table (Dridex)
 Dynamic Resolving (Shifu)
 Imports on Stack / Heap (PIVY, Cryptowall)
 XORed Imports (Qadars)

Covered by ApiScout [1]



[1] https://github.com/danielplohmann/apiscout

Occurrence Frequency of Individual WinAPI Functions

Occurrence frequency per Windows API function



- Discounting .NET and API-obfuscated families
 - Only 3 API functions > 90% (CloseHandle, Sleep, WriteFile)
 - Only 48 API functions > 50%
 - API function at position 150 appears in 21.73%
 - 4,392 (92.52%) of API functions <= 10%
- API compositions are highly specific per family
 - Good for (identification) tools like
 - ImpHash [1]
 - ImpFuzzy [2]
 - ApiVectors!



[1] https://www.fireeye.com/blog/threat-research/2014/01/tracking-malware-import-hashing.html

15 [2] http://blog.jpcert.or.jp/2017/03/malware-clustering-using-impfuzzy-and-network-analysis---impfuzzy-for-neo4j-.htm

Semantic Context for Windows API Functions

- Define: API Context Groups
 - Manually labelled ~4.500 APIs, primary (12) and secondary class (115)

GUI	1392
System	636
Execution	590
String	458
Network	387
FileSystem	352
Device	170
Crypto	131
Other	127
Memory	118
Registry	80
Time	44

Kudos to Quoscient.io for their contributions! (Patrick Ventuzelo, Lukas Bernhard)



16

WinAPI Reference Vector

- Define: API Context Groups
 - Reduce this set to 1024 WinAPIs (~80% hierarchy, ~20% based on domain knowledge)

GUI	1392		Execution	229
System	636		Network	192
Execution	590	$ \langle \ \rangle $	System	150
String	458		FileSystem	114
Network	387		Memory	68
FileSystem	352		Device	66
Device	170		String	52
Crypto	131		Crypto	48
Other	127	\bigvee	Registry	32
Memory	118	l' > 4	GUI	27
Registry	80		Other	24
Time	44] →	Time	22



WinAPI Reference Vector

- Define: API Context Groups
 - Reduce this set to 1024 WinAPIs (~80% hierarchy, ~20% based on domain knowledge)
 - Vector yields 90% coverage (mean) for APIs found by ApiScout for ~600 malware families





ApiScout WinAPI Reference Vector

Visualize Vectors:

Hilbert Curve to ensure neighboring of contexts







ApiScout WinAPI Reference Vector

Visualize Vectors:

Hilbert Curve to ensure neighboring of contexts







ApiScout WinAPI Reference Vector

Some Examples





Vector Construction





?

Α

В

AUB

22

Comparison of ApiVectors

- Example Vectors
 - Base64-like encoding (Run-Length compressed) 4-172 bytes long

A42gA28KA13 CAAMA16BABAAJAECAxMAACkAAQUA7CJBCgAgUBA3 kQCBAHJSRjU^q-*}_pb__N,__^? A42gA28KA13 CAAMA16BABAAJAEAAxMAACkAAQUA7CJBCgAAUBA3 kQCBAHJSRjU^q-*}_pL__N,._^? A41BA29CA4IA9gCA9gA8Q BAAJAEAABMA3 gAAQA8 QJRCgAgUBAAHkQARCDIADDBGAqQAgCcGOIOp,f?

TeslaCrypt 2.2, 3.0, 4.2





Evaluation of Matching Performance





24

Evaluation of Matching Performance



- General Challenges to API-based similarity analysis
 - Packers
 - .NET / scripts
 - Statically linked code (MSVCRT, Delphi, Go, ...)



25

Windows API Usage Recovery & Analysis

How to operationalize this?

- ApiScout available on GitHub [1]
- Projects using ApiScout:
 - Angad [2] by Ankur Tyagi, presented @ BsidesZurich [3]
 - Master of Clusters by Andrea Garavaglia to be presented @ MISP Summit / hack.lu [4]
 - Malpedia :)
- Ideally used as post-processing for sandboxing
 - Cuckoo plugin?



Malpedia Usage Example: Code-based Similarity Analysis



Code-based Similarity Analysis

Motivation

Code Similarity Analysis!!!

- Identify (3rd party) shared library code: automated annotation / exclusion from analysis scope
- Isolate code that is immanent to a given code base / author

Related Work:

- Kam1n0 [1] by Stephen Ding et al.
- FunctionSimSearch [2] by Thomas Dullien et al.
- CosaNostra / MalTindex [3] by Joxean Koret
- More...



Code-based Similarity Analysis

Overview

Tool: SMDA [2]

- Work in progress, already silently released on GitHub [2]
- Built on top of Capstone [1]
- "SMDA is a minimalist recursive disassembler library that is optimized for accurate Control Flow Graph (CFG) recovery from memory dumps."

Tool: MCRIT

- "MinHash-based Code Relationship Identification Toolkit"
- Work in progress, will be released later this year



Code-based Similarity Analysis SMDA

- Malpedia uses memory dumps as normalization
 - Available tools (IDA, Binary Ninja, radare2, ...) not suited / optimized for this

```
/home/pnx/smda $ python3 analyze.py citadel_dump_0x00140000
        -o citadel_disasm.json
        now analyzing /home/pnx/citadel_dump_0x02390000
```

2018-09-24 11:30:57,299: smda.common.ApiResolver - loaded 57315 exports from 134 DLLs (Windows XP Professional).

-> 3.86s | 873 Func (status: ok)

/home/pnx/smda \$ cat citadel_disasm.json

Example results:

Method	FNC	BB	INS	TPR
Manual Reference:	793	10,264	52,121	99.99%
SMDA vs. dump:	873	10,686	53,324	98.36%
IDA vs. dump:	658	9,071	46,266	85.62%
IDA vs. clean unpacked:	788	10,209	51,794	99.37%

"architecture": "intel". "base addr": 1310720. "bitness": 32. "execution time": 3.86346, "filename": "citadel_dump_0x00140000", "message": "Analysis finished regularly.", "sha256"; "2c4166c81c31fd5de7fe205ea12c74d8bb394005337fa844afc23e4707a3d42f". "status": "ok". "summary": { "num_api_calls": 1146, "num basic blocks": 10686. "num disassembly errors": 56, "num_function_calls": 3979 "num functions": 873. "num instructions": 53324 "num leaf functions": 114. "num recursive functions": 7 "timestamp": "2018-09-24T09-31-01", "version": "1.0.1", "xcfg": { "1318184": { "apirefs": {}. "blockrefs": {}, "blocks": { "1318184": [1318184. "55", "push". "ebp"



Code-based Similarity Analysis

MinHash 101

- MinHashing
 - "Min-wise independent permutations" Locality Sensitive Hashing (LSH) scheme [1]
 - Fast estimation of set similarity (approximation of Jaccard similarity coefficient)

Use cases:

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



Code-based Similarity Analysis

MinHash 101

- MinHash procedure:
 - Extract a range of descriptive features ("shingles") for each object
 - Hash them n times (with each hash function seeded differently)
 - Select the minimum hash value for each of the n groups
 - The resulting sequence of n values is considered as the function's fingerprint

Matching fingerprints:

Given two fingerprints, count the number of equal fields at same positions

Various optimizations:

Single-hash XORing, Banding or n-key sorting, b-bit representation, ...



Code-based Similarity Analysis MCRIT

Simplified example with a hash function that maps to a single output byte (0-255)



Code-based Similarity Analysis

- Small test data set (in-memory):
 - 50 samples, 40 families
 - 26,097 functions with 20,611 indexable (greater or equal to 10 instructions or 3 basic blocks)

Application of MCRIT

- All function pairs: 20,611 * 20,610 / 2 = 212,396,355
- Filter candidates down to 35,651 pairs (using banding)
- This results in 19,732 matches above threshold (0.7)
- Indexing + Matching takes ~2min on this laptop (i5, 8GB RAM).
- Comparison: BinDiff
 - All samples vs. each other (50*49/2 = 1,225 pairs)
 - Runtime: ~60min

BinDiff Threshold	0.90	0.99
BinDiff Matches	12,035	8,263
MCRIT Threshold	0.70	0.85
MCRIT Matches	19,732	11,648
MCRIT TPs	9,350	7,968
MCRIT TPR	0.7769	0.9643
MCRIT FPs (?)	3,515	766

Preliminary Results!



Code-based Similarity Analysis MCRIT

- Malpedia data set (mongodb):
 - 2,403 samples, 773 families
 - 1,927,361 functions with 1,233,321 indexable (greater or equal to 10 instructions or 3 basic blocks)

Application of MCRIT

- All function pairs: 1,233,321 * 1,233,320 / 2 = 760,539,727,860
- Filter candidates down to 63,694,525 pairs (banding)
- This results in 29,596,574 matches above threshold (0.7)

Runtime

35

- Indexing: 13,902 sec (03:51:42h) 138,64 FNs/sec
- Candidate Identification: 6,380 sec (01:46:20h)
- Matching: 31,840 sec (08:50:40h) 1666,52 Pairs/sec
- Basically, on this laptop.

Code-based Similarity Analysis

- Next steps
 - Improve matching quality
 - Tweak / verify against multiple ground truth data sets
 - Goodware / libraries with different compilers available
 - Make it usable
 - REST API
 - integrations with other analysis tools (IDA, Binja, r2, ...?)
 - Extensive evaluation on Malpedia data set
- Hosted service along Malpedia?



Summary



Summary Code Cartographer's Diary



The Malpedia Vision: A curated, free, high-quality malware corpus for research

- Vetted + community-driven
- Want Access?
 - Talk to me (Know Met Trust (KMT) -> ensures K&M already)
 - Get an invite by another existing member that can vouch for you
 - Procedure can be potentially accelerated based on your background (GOV/LEA, ...)
- Windows API Usage Recovery & Analysis
 - ApiScout: Convenient & reliable WinAPI usage recovery from memory dumps
 - ApiVectors: Compact representation, decent matching performance
- Code-based Similarity Analysis
 - SMDA: Recursive disassembler (FOSS) optimized for memory dumps
 - MCRIT: Code-based similarity analysis has huge potential



Thank You for Your Attention!

Daniel Plohmann daniel.plohmann@fkie.fraunhofer.de





