



Contents lists available at ScienceDirect

Computers & Security

journal homepage: www.elsevier.com/locate/cose

A Systematical and longitudinal study of evasive behaviors in windows malware

Nicola Galloro^a, Mario Polino^{a,*}, Michele Carminati^a, Andrea Continella^b, Stefano Zanero^a

^a Politecnico di Milano, Dipartimento di Elettronica, Informazione e Bioingegneria (DEIB), Via Ponzio 34/5, Milano, 20133, Italy

^b University of Twente, Faculty of Electrical Engineering, Mathematics and Computer Science (EEMCS), Zilverling building, Drienerolaan 5, 7522 NB, Enschede, The Netherlands

ARTICLE INFO

Article history:

Received 10 March 2021

Revised 4 September 2021

Accepted 14 November 2021

Available online 21 November 2021

ABSTRACT

Malware is one of the prevalent security threats. Sandboxes and, more generally, instrumented environments play a crucial role in dynamically analyzing malware samples, providing key threat intelligence results and critical information to update detection mechanisms.

In this paper, we study the evasive behaviors employed by malware authors to hide the malicious activity of samples and hinder security analysis. First, we collect and systematize 92 evasive techniques leveraged by Windows malware to detect and thwart instrumented environments (e.g., debuggers and virtual machines). Then, we implement a framework for evasion analysis of x86 binaries and analyze 45,375 malware samples observed in the wild between 2010 and 2019; we compare this analysis against popular, legitimate Windows programs to study the intrinsic characteristics of such evasive behaviors.

Based on the results of our experiments, we present statistics about the adoption of evasive techniques and their evolution over time. We show that over the past 10 years, the prevalence of evasive malware samples had a slight increase (12%). Moreover, the employed techniques shifted significantly over time. We also identify techniques that are specific to malware, as opposed to being employed by both malicious and legitimate software. Finally, we study how the security community reacts to the deployment of new evasive techniques. Overall, our results empirically address open research questions and provide insights and directions for future research.

© 2021 Elsevier Ltd. All rights reserved.

1. Introduction

Malware is one of the most dangerous and prevalent Internet threats. Financially-motivated threats, such as ransomware (Continella et al., 2016; Kharraz et al., 2015), banking trojans (Continella et al., 2017), and cryptominers (Holz et al., 2020; Kharraz et al., 2019; Konoth et al., 2018), have dominated the malicious landscape in the past years (Labs, 2017; McAfee.com, 2018). Given the complexity of applying static analysis to modern malware samples, especially when heavily obfuscated (Ugarte-Pedrero et al., 2015), malware analysts often rely on dynamic analysis to reverse-engineer new samples, providing key threat intelligence results and critical information to update detection mechanisms. In fact, dynamic analysis is starting to be also employed in malware detectors, some of which leverage emulators to execute

suspicious binaries in an instrumented environment (Quarta et al., 2018).

Previous literature studied how malicious samples employ techniques to recognize when they are executed in a controlled environment and subsequently hide or alter their malicious behavior. The detection of the analysis environments is performed through *fingerprinting*: Malware samples look for specific artifacts left by the analysis components (or agents). If any such artifact is identified, the malicious activity is not executed (or executed differently) (Lindorfer et al., 2011; Miramirkhani et al., 2017; Paleari et al., 2009; Raffetseder et al., 2007; Rutkowska, 2004). Different solutions have been proposed to address this problem. Multiple path exploration approaches force the execution of the whole program to discover hidden behaviors but suffer from the well-known path explosion problem (Comparetti et al., 2010). An alternative is to detect inconsistencies between the execution of a sample in a sandbox and the execution of the same sample without the instrumentation layer. This approach has been tested by Kirat et al. (2014), who measured the difference between bare-metal and sandbox executions to detect “paths” related to eva-

* Corresponding author.

E-mail address: mario.polino@polimi.it (M. Polino).

sion techniques. However, this method does not intercept low-level system interactions, such as system calls and x86 instruction, because the bare-metal setup can only reason about raw disk content. A more recent solution is to execute the malicious samples inside a Dynamic Binary Instrumentation (DBI) framework that hides the presence of the analysis environment, identifying anti-instrumentation techniques (D'Elia et al., 2019; 2020; Polino et al., 2017).

This paper provides the first comprehensive, systematic, and longitudinal study of the evasive behaviors adopted by Windows malware. Our research attempts to answer important questions that are currently still open: What are the prevalent evasion techniques adopted by malware samples to hinder security analysis? How have such techniques evolved in the past years? How does the security community react to the adoption of a new evasive technique? Are evasive techniques also adopted by legitimate programs? To answer these questions, we first collect and systematize 92 evasive techniques adopted by Windows malware. To the best of our knowledge, we provide the most extensive taxonomy of evasive techniques: The previous most complete surveys (Afianian et al., 2019; Bulazel and Yener, 2017) cover 41 techniques. Then, we develop a fine-grained framework to identify and collect information about the evasive behaviors present in Windows programs. Built on top of Arancino (Polino et al., 2017), our framework leverages the capabilities of a DBI tool, Intel Pin, to instrument the target executable and identify when the program under analysis attempts to collect artifacts to evade the analysis itself.

Using our framework and taxonomy, we study the evasive behaviors adopted by 45,375 malware samples observed in the wild between 2010 and 2019. We identify the most common evasive techniques and correlate them with the different malware families, providing insights about different malware campaigns. To further investigate the characteristics of such evasive techniques, we test their adoption in popular, legitimate Windows programs. Our longitudinal analysis indicates that the amount of evasive malware samples has not increased significantly in the past 10 years. However, the techniques adopted by these evasive samples varied over time.

In summary, we make the following contributions:

- We systematically study evasive behaviors employed by Windows malware by collecting 92 techniques used to evade instrumentation environments (e.g., Debuggers, Virtual Machines)—the largest collection so far—and providing a taxonomy to classify such techniques according to their semantics and characteristics.
- We analyze 45,375 malware samples observed in the wild and 949 popular, legitimate applications. We study the correlation between malware families and the adopted evasive techniques, and we analyze both malicious and legitimate programs to compare their evasive behaviors.
- We perform a longitudinal study of the adoption of evasive techniques in malware families between 2010 and 2019, and we provide statistics and insights about the evolution of evasive behaviors. Additionally, we empirically estimate how the community reacted to and influenced the adoption of evasive techniques during the same period. For instance, we show that the adopted techniques varied in the past years, while the number of evasive samples is rather stable.

2. Related work

The analysis of the evasive techniques by which malware determines the presence or absence of an analysis system is a wide research topic for which we refer the reader to (Afianian et al., 2019; Bulazel and Yener, 2017) that provide comprehensive surveys

and taxonomies of 41 existing anti-analysis techniques. Unlike our work, which quantifies the adoption of 92 evasive techniques from the dynamic execution of actual malware, they put in place a qualitative analysis only. In this section, we describe the related works that investigate the evasive behavior in a large-scale study of malware found in the wild. Table 1 provides a comparison of the recent works in the area, according to the type of analysis they conducted. Regarding the analysis of evasive techniques used by legitimate non-malicious software vendors in order to prevent analysis or piracy of their software, we refer the reader to (Blackthorne et al., 2016; Collberg and Thomborson, 2002; Gagnon et al., 2007).

Xu Chen et al. (2008) proposed a first taxonomy of malware anti-analysis techniques by comparing execution traces from different environments. They found that 2.7% of 6,222 samples from 2006–2007 exhibited less behavior when running in a VM versus a bare-metal system.

Branco et al. (2012) and Barbosa and Branco (2014) introduced various static detection methods for anti-debugging and anti-VM techniques and run analysis over 4 and 12 million samples, respectively, to show the state of evasion techniques in use. In their study, anti-analysis operations were detected only through static analysis, while dynamic behavior was not considered. The main problem with static analysis resides in packing and obfuscation techniques, which are widely used by malware authors precisely to avoid it. Indeed, since static analysis is based on parsing and comparing executable byte sequences, obfuscated samples are opaque to it. Additionally, there is no guarantee that any matched technique is actually executed. On the contrary, with dynamic analysis, we can precisely identify evasive techniques while they are executed, even when packing and obfuscation are applied. In summary, dynamic analysis is more precise in this context. Compared to static analysis, our approach—and dynamic analysis in general—is resilient to packing due to the dynamic nature of the DBI analysis that allows a profound control over the instrumented binary.

Chen et al. (2016) surveyed anti-VM and anti-debugging techniques in 17,283 “generic” and targeted (APT) malware samples collected between 2009–2014. The authors found that these techniques became more prevalent over the surveyed timeframe and that their presence was negatively correlated with antivirus detection. They also observed that APTs generally do not use as many anti-analysis techniques as generic malware and that they have decreased their use over time. Although their study is similar to the present one, since they provide a longitudinal study of the trends in anti-analysis techniques analyzing malware samples, there are also many differences. First, in their study, anti-analysis techniques were detected only through static analysis, while dynamic behavior was not considered. Therefore, it shares the same limitations of (Barbosa and Branco, 2014; Branco et al., 2012) presented before. Moreover, they limit their analysis to only six malware families without clarifying in detail the techniques executed by each of them. Consequently, their results are biased by the restricted number of families considered. In the present work, instead, we strive to select malware samples equally distributed in time.

Polino et al. (2017) found that 15.6% of 7,006 malware samples collected from VirusTotal between October 2016 and February 2017 exhibited anti-instrumentation behavior (not necessarily solely anti-DBI). Similarly to the present work, the authors leverage the capabilities of a DBI tool, Intel Pin, to instrument the target executable and identify when the program under analysis attempts to collect artifacts to evade the analysis itself. However, they do not consider Anti-VM and Anti-Debugging behaviors.

Miramirkhani et al. (2017) assesses the threat of malware sandbox evasion strategies that leverages artifacts indicative of the “wear and tear” and “aging” of a system to identify artificial environments and demonstrate its effectiveness against existing sandboxes used by malware analysis services. Commonly to this work,

Table 1

Comparison of large-scale studies of evasive malware techniques: ● dynamic analysis, ○ static analysis, ✓/✗ if the work provides or not, respectively, the considered analysis.

	Year	Taxonomy	Dynamic-Static Analysis	Longit. Analysis	Malware vs Goodware	# Evasive Techniques	# Samples (# Families)
Chen et al. Xu Chen et al. (2008)	2006	✓	●	✗	✗	8	6,222 (✗)
Lindorfer et al. Lindorfer et al. (2011)	2011	✓	●	✗	✗	14	1,500 (175)
Branco et al. Branco et al. (2012)	2012	✗	○	✗	✗	51	4M (✗)
Barbosa et al. Barbosa and Branco (2014)	2014	✗	○	✗	✗	51	12M(✗)
Polino et al. Polino et al. (2017)	2017	✗	●	✗	✗	24	7006 (✗)
Miramirkhani et al. Miramirkhani et al. (2017)	2017	✓	●	✗	✗	44	✗
Oyama Oyama (2018b)	2018	✗	●	✗	✗	28	8243 (3044)
Afianian et al. Afianian et al. (2019)	2019	✓	✗	✗	✗	46	✗
D'Elia et al. D'Elia et al. (2020)	2020	✓	●	✗	✗	45	1000 (✗)
Chen et al. Chen et al. (2016)	2021	✓	○	2009-2014	✗	41	17,283 (6)
This work	2021	✓	●	2010-2019	✓	92	45,375(2867)

Table 2

Top 20 families. # refers to the number of samples, **Date** refers to the first appearance in VirusTotal of a sample belonging to the specified family, # **Ev.** is the number of evasive samples exposing at least a technique, # **Packed** is the number of samples that are packed according to PEID, # **Tech.** is the maximum number of techniques detected in a sample of the related family.

Family	#	Date	# Ev.	# Packed	# Tech.
zbot	2412	01/2010	1564 (65%)	601 (25%)	19
vtflooder	1016	06/2014	991 (98%)	1 (0%)	15
installex	859	01/2012	850 (99%)	2 (0%)	14
firseria	835	08/2011	801 (96%)	91 (11%)	11
multiplug	807	03/2012	782 (97%)	93 (12%)	14
fareit	750	08/2011	624 (83%)	506 (67%)	21
domaiq	743	02/2013	741 (99%)	67 (9%)	7
high	652	01/2010	652 (53%)	552 (45%)	19
vobfus	570	01/2010	227 (40%)	444 (78%)	15
delf	558	01/2010	490 (88%)	462 (83%)	17
flystudio	550	01/2010	474 (86%)	474 (86%)	20
dealply	485	12/2010	424 (87%)	433 (89%)	15
neshta	482	06/2011	470 (98%)	450 (93%)	17
installcore	474	06/2014	356 (75%)	330 (70%)	10
gamarue	473	10/2011	348 (74%)	273 (57%)	11
xorist	447	10/2010	48 (11%)	15 (3%)	12
loadmoney	440	12/2012	367 (83%)	62 (14%)	14
autoit	436	01/2010	420 (96%)	28 (6%)	18
bladabindi	404	02/2011	83 (21%)	343 (85%)	15
virlock	402	11/2014	272 (68%)	1 (0%)	15

as presented in [Tables 3 and 4](#), we check the “wear-and-tear” state of the machine. However, while we focus on studying the evasion techniques currently exploited by malware, this work, instead, proposes novel evasive techniques against the existing sandbox mechanism.

[Oyama \(2018b\)](#) analyzes a dataset of dynamic malware analysis results and examines trends in the Anti-Analysis operations executed by malware samples collected in 2016. Among 8243 malware samples, 856 (10.4%) samples executed at least one type of the 28 Anti-Analysis operations investigated.

[D'Elia et al. \(2020\)](#) present a DBI-based prototype that, similarly to [Arancino \(Polino et al., 2017\)](#) and our work, offers a stealthy execution environment based on an observe-check-replace design to intercept evasive attempts and hide imperfections. They demonstrate its effectiveness against a set of 45 highly evasive samples.

In contrast to the presented studies, in this work, we perform a longitudinal study of the adoption of evasive techniques by dynamically executing malware to provide quantitative results and insights into their evolution and adoption. Moreover, we explore the adoption of evasive behaviors in legitimate programs to investigate whether evasive techniques are intrinsically evasive or also used for legitimate purposes. Finally, we analyze a higher number of samples and techniques.

3. Taxonomy of evasive techniques

An evasive technique is a specific activity performed by a malware sample to detect a dynamic analysis system (and subsequently thwart it). In particular, we focus on any activity, under the form of system calls, Windows API functions, memory accesses, or x86 instructions, that leaks the presence (i.e., artifacts) of analysis environments. We study evasive techniques targeting virtual machines and sandboxes (i.e., Anti-VM techniques), as well as debuggers and DBI implementations (i.e., Anti-Debugging techniques). Note that some techniques can overlap this broad categorization because of common artifacts leaked by both classes of environments.

After gathering and analyzing 92 evasive techniques found in the literature, we propose to classify them according to the semantics of the action(s) performed on the OS. Proceeding thus, starting from systematization previously presented ([Afianian et al., 2019](#); [Chen et al., 2016](#); [D'Elia et al., 2019](#); [Lindorfer et al., 2011](#); [Miramirkhani et al., 2017](#); [Xu Chen et al., 2008](#)), we develop a taxonomy consisting of 16 semantic-equivalent groups. In the following, we list the groups, providing a few examples of techniques that belong to each. For further details on each technique, we refer the reader to [Appendix A](#).

Memory Fingerprinting. These techniques access memory regions of the running process to detect debuggers. In fact, debuggers modify several system variables associated with their target process. For instance, the PEB, a structure present in each Windows process, contains variables designed to indicate the presence of a debugger.

Exception Handling. With this type of approach, a malware sample throws exceptions that a debugger handles for debugging purposes. The malicious sample sets up a handling routine for the exception and, if the routine is never executed, it means that the debugger caught the exception. A less trivial example is `NtClose(INVALID_HANDLE)`: Closing an invalid handle causes an exception only if there is a debugger attached to the process.

CPU Fingerprinting. This category contains all the Anti-VM techniques that fingerprint the CPU, detecting differences between physical CPUs and virtualized ones. This behavior can be intended, or it can be a side-effect due to incorrect execution of sensitive instructions by the hypervisor. An example is `CPUID`, an instruction that retrieves information and features from the CPU and returns the presence of a hypervisor and its brand if present ([CheckPoint, 2020](#)).

Table Descriptors. This category contains x86 instructions used to retrieve the addresses of OS Table Descriptors. These addresses frequently change if a hypervisor is present. Thus they have been used to detect virtualized environments. Nevertheless, these tech-

Table 3

Techniques **D/V**: ● Anti-Debugging, ● Anti-VM, ● Both. ◇ Used by Goodware. **BB** is the number of basic blocks, and **I** is the number of instructions, **C** is the number function calls, **ASM** indicates whether implementing the given technique requires writing assembly code.

	Technique	#Families (%)	#Goodware (%)	#Malware (%)	D/V	G.	First	Last	BB	I	C	ASM	
Memory Fingerprinting	Kulchytskyy and Kukoba (2019); Microsoft (2020); Xu Chen et al. (2008)	PEB→IsDebugged	162 (5.65)	72.0 (7.58)	982 (2.16)	●	◇	01/10	09/19	1	9	0	Y
		KUSER_SHARED→KdDebuggerEnabled	23 (0.8)	0 (0)	102 (0.22)	●		01/11	08/19	1	8	0	Y
		PEB→NtGlobalFlag	46 (1.6)	0 (0)	96 (0.21)	●		02/10	03/19	1	9	0	Y
		PEB→Heap→Flags	2 (0.07)	0 (0)	3 (0.01)	●		05/17	02/19	4	16	0	Y
		PEB→Heap→ForceFlags	2 (0.11)	0 (0)	3 (0.01)	●		09/16	05/17	1	10	0	Y
Exception Handling	Ferrie (2011)	SetUnhandledExceptionFilter	1817 (63.38)	475.0 (50.0)	15651 (34.49)	●	◇	01/10	09/19	8	27	3	Y
		NtClose(INVALID_HANDLE)	606 (21.14)	19 (2.0)	3105 (6.84)	●	◇	01/10	09/19	14	47	6	N
		OutputDebugString	234 (8.16)	26.0 (2.74)	794 (1.75)	●	◇	01/10	08/19	14	55	6	N
		POPF/D - TRAP FLAG	58 (2.02)	0 (0)	180 (0.4)	●		01/10	06/19	9	38	4	Y
		CTRL Exception Handling	0 (0)	0 (0)	0 (0)	●		Never	Never	10	34	4	N
CPU Fingerprinting	Cyberbit (2007); Sikorski and Honig (2012)	CPUID(EAX=0x00000001)	1333 (46.49)	320.0 (33.68)	10070 (22.19)	●	◇	01/10	09/19	1	11	0	Y
		IN	234 (8.16)	0 (0)	869 (1.92)	●		01/10	09/19	12	51	5	Y
		CPUID(EAX=0x40000000)	31 (1.08)	0 (0)	119 (0.26)	●		12/10	08/19	23	171	7	Y
		STR	50 (1.74)	0 (0)	122 (0.27)	●		01/10	09/18	5	18	2	Y
		SMSW	9 (0.31)	0 (0)	11 (0.02)	●		02/10	06/15	5	18	2	Y
Table Descr.	Sikorski and Honig (2012)	SLDT	37 (1.29)	0 (0)	176 (0.39)	●		01/10	08/19	1	5	0	Y
		SGDT	16 (0.56)	0 (0)	126 (0.28)	●		03/11	08/19	6	28	1	Y
Traps	Ferrie (2011); Instinct (2017); Kulchytskyy and Kukoba (2019); Sikorski and Honig (2012); vpc (2017)	SIDT	46 (1.6)	0 (0)	162 (0.36)	●		01/10	04/17	6	22	1	Y
		INT 3	172 (6.0)	0 (0)	747 (1.65)	●		01/10	08/19	11	38	4	Y
		VPCEXT	105 (3.66)	0 (0)	362 (0.8)	●		01/10	09/19	14	52	6	Y
		POP SS	92 (3.21)	0 (0)	289 (0.64)	●		01/10	05/19	1	14	0	Y
		INT 1	74 (2.58)	0 (0)	141 (0.31)	●		02/10	08/19	11	37	4	Y
		ICEBP	27 (0.94)	0 (0)	61 (0.13)	●		01/10	09/19	11	37	4	Y
		INT 2D	10 (0.35)	0 (0)	39 (0.09)	●		03/10	05/19	11	38	4	Y

(continued on next page)

Table 3 (continued)

Timing	McAfee (2019); Polino et al. (2017); Sikorski and Honig (2012)	GetTickCount	1464 (51.06)	198.0 (20.84)	11029 (24.31)	●	◇	01/10	09/19	9	34	3	N		
		RDTSC/D	1398 (48.76)	168.0 (17.7)	9518 (20.98)	●	◇	01/10	09/19	31	231	3	Y		
		QueryPerformanceCounter	994 (34.67)	140.0 (14.74)	6038 (13.31)	●	◇	01/10	09/19	41	243	6	N		
		GetLocalTime	587 (20.47)	75.0 (7.89)	2707 (5.97)	●	◇	01/10	09/19	9	41	3	N		
		timeGetTime	217 (7.57)	32 (3.37)	805 (1.77)	●	◇	01/10	09/19	3	10	1	N		
		NtQuerySystemTime	174 (6.07)	20.0 (2.11)	667 (1.47)	●	◇	01/10	09/19	9	50	3	N		
		GetSystemTimes	22 (0.77)	0 (0)	37 (0.08)	●		07/13	09/19	9	43	3	N		
		KUSER_SHARED→SystemTime	13 (0.45)	0 (0)	27 (0.06)	●		01/13	09/19	6	61	1	Y		
		KUSER_SHARED→InterruptTime	11 (0.38)	0 (0)	19 (0.04)	●		04/10	06/19	6	61	1	Y		
		timeGetSystemTime	3 (0.1)	0 (0)	6 (0.01)	●		07/10	12/17	9	43	3	N		
		KUSER_SHARED→TickCountQuad	12 (0.42)	0 (0)	17 (0.04)	●		09/10	05/17	6	30	1	Y		
		NtGetTickCount	0 (0)	0 (0)	0 (0)	●		Never	Never	13	44	5	N		
		QueryInterruptTime	0 (0)	0 (0)	0 (0)	●		Never	Never	15	67	5	N		
		NtQueryPerformanceCounter	0 (0)	0 (0)	0 (0)	●		Never	Never	36	230	7	N		
		QueryUnbiasedInterruptTimePrecise	0 (0)	0 (0)	0 (0)	●		Never	Never	15	67	5	N		
		QueryInterruptTimePrecise	0 (0)	0 (0)	0 (0)	●		Never	Never	15	67	5	N		
		GetTickCount64	0 (0)	0 (0)	0 (0)	●		Never	Never	15	61	5	N		
		QueryUnbiasedInterruptTime	0 (0)	0 (0)	0 (0)	●		Never	Never	15	67	5	N		
		Stalling	Kolbitsch et al. (2011); Oyama (2018a)	waitForSingleObject/Ex	1291 (45.03)	242.0 (25.47)	9117 (20.09)	●	◇	01/10	09/19	9	36	4	N
				Sleep/SleepEx	1249 (43.56)	41.0 (4.32)	9135 (20.13)	●	◇	01/10	09/19	3	10	1	N
SetTimer	842 (29.37)			100.0 (10.53)	5686 (12.53)	●	◇	01/10	09/19	14	45	4	N		
SetWaitableTimer/Ex	75 (2.62)			1.0 (0.11)	195 (0.43)	●		01/10	09/19	7	46	3	N		
CreateTimerQueueTimer	23 (0.8)			2.0 (0.21)	92 (0.2)	●		04/11	02/19	20	64	6	N		
NtDelayExecution	22 (0.77)			37.0 (3.89)	130 (0.29)	●	◇	07/10	05/19	7	39	3	N		
timeSetEvent	33 (1.15)			6.0 (0.63)	84 (0.19)	●		02/10	08/19	5	22	1	N		
icmpSendEcho/2/Ex	6 (0.21)			1.0 (0.11)	42 (0.09)	●		11/13	07/19	16	54	4	N		
GetCursorPos	754 (26.3)			99.0 (10.42)	4157 (9.16)	●	◇	01/10	09/19	10	32	3	N		
Registry	Instinct (2019)			GetLastInputInfo	49 (1.71)	1.0 (0.11)	305 (0.67)	●		04/10	09/19	6	26	3	N
		NtOpenKey/Ex	193 (10.55)	12.0 (1.26)	781 (3.13)	●	◇	01/13	09/19	24	178	5	N		
		NtEnumerateKey	215 (7.5)	2.0 (0.21)	790 (1.74)	●		01/10	09/19	69	334	18	N		
		NtQueryValueKey	93 (3.24)	9.0 (0.95)	304 (0.67)	●	◇	03/10	09/19	16	70	5	N		
		NtEnumerateValueKey	6 (0.21)	0 (0)	58 (0.13)	●		07/13	10/18	19	78	3	N		
System Environment	Cyberbit (2007); Lindorfer et al. (2011); McAfee (2019); Project (2017a); Sikorski and Honig (2012)	NtQuerySystemInformation(PHYSICAL_MEMORY_INFO)	1342 (46.81)	157.0 (16.53)	10651 (23.47)	●	◇	01/10	09/19	15	53	4	N		
		FindWindow	520 (18.14)	51.0 (5.37)	2245 (4.95)	●	◇	01/10	09/19	10	34	3	N		
		GetComputerName	403 (14.06)	18.0 (1.9)	2041 (4.5)	●	◇	01/10	09/19	15	50	2	N		
		IsDebuggerPresent	304 (10.6)	17.0 (1.79)	1516 (3.34)	●		01/10	09/19	6	14	1	N		
		GetAdaptersInfo	165 (5.76)	2.0 (0.21)	1328 (2.93)	●		01/10	09/19	54	136	13	N		
		GetDiskFreeSpace/Ex	184 (6.42)	47.0 (4.95)	528 (1.16)	●	◇	01/10	08/19	7	30	1	N		
		CheckRemoteDebuggerPresent	124 (4.33)	1 (0.11)	432 (0.95)	●		01/10	09/19	9	24	2	N		
		GetAdaptersAddresses	58 (2.02)	4 (0.42)	293 (0.65)	●	◇	02/11	09/19	46	135	13	N		
		GlobalMemoryStatusEx	78 (2.72)	26.0 (2.74)	250 (0.55)	●	◇	11/10	09/19	5	32	1	N		
		NtQuerySystemInformation(0x23)	28 (0.98)	0 (0)	49 (0.11)	●		03/10	09/19	10	33	3	N		
		DeviceIoControl	2 (0.07)	2 (0.21)	7 (0.02)	●		07/13	11/17	14	62	4	N		
		PEB→NumProcessors	7 (0.24)	0 (0)	11 (0.02)	●		03/12	11/17	1	13	0	Y		
		NtQueryObject after NtCreateDebugObject	0 (0)	0 (0)	0 (0)	●		Never	Never	21	93	6	N		
		ExecQuery	0 (0)	0 (0)	0 (0)	●		Never	Never	25	114	10	N		
		WMI	FireEye (2016)												

Table 4

Techniques **D/V**: Anti-Debugging, Anti-VM, Both. Used by Goodware. **BB** is the number of basic blocks, and **I** is the number of instructions, **C** is the number function calls, **ASM** indicates whether implementing the given technique requires writing assembly code.

	Technique	#Families (%)	#Goodware (%)	#Malware (%)	D/V	G.	First	Last	BB	I	C	ASM	
Process Environment	Ferrie (2011); Kulchytsky and Kukoba (2019); Sikorski and Honig (2012)	NtSetContextThread(CONTEXT_DEBUG_REGISTERS)	257 (8.96)	0 (0)	1617 (3.56)			01/10	09/19	13	40	2	N
		NtQueryInformationProcess(0x07)	280 (9.77)	4.0 (0.42)	1028 (2.27)			01/10	09/19	13	42	4	N
		NtSetInformationThread(0x11)	122 (4.26)	1.0 (0.11)	350 (0.77)			01/10	07/19	14	46	7	N
		NtQueryInformationProcess(0x1e)	96 (3.35)	0 (0)	303 (0.67)			01/10	08/19	15	56	6	N
		NtQueryInformationProcess(0x1f)	34 (1.19)	0 (0)	115 (0.25)			10/12	08/19	13	41	4	N
		NtGetContextThread(CONTEXT_DEBUG_REGISTERS)	28 (0.98)	0 (0)	46 (0.1)			01/10	03/18	14	42	2	N
File System	Instinct (2019)	DebugActiveProcess on Parent	2 (0.11)	0 (0)	2 (0.01)			10/14	02/16	54	154	21	N
		NtCreateThreadEx(HIDE_FROM_DEBUGGER)	0 (0)	0 (0)	0 (0)			Never	Never	16	50	6	N
		NtOpenFile	613 (21.38)	6.0 (0.63)	3268 (7.2)			01/10	09/19	7	40	2	N
		NtQueryAttributesFile	536 (18.7)	14.0 (1.48)	2616 (5.77)			01/10	09/19	7	34	2	N
		NtCreateFile	389 (13.57)	6.0 (0.63)	1696 (3.74)			01/10	09/19	9	49	2	N
		NtQueryDirectoryFileEx	0 (0)	0 (0)	0 (0)			Never	Never	81	330	23	N
List Proc.	Instinct (2019)	NtQSI (SYSTEM_PROCESS_INFO)	649 (22.64)	29.0 (3.06)	4054 (8.93)			01/10	09/19	47	189	22	N
		EnumProcesses	122 (4.26)	11.0 (1.16)	393 (0.87)			02/10	12/18	58	201	16	N
List Services	Instinct (2019)	GetModuleBaseName	0 (0)	0 (0)	0 (0)			Never	Never	72	254	25	N
		OpenSCManager	222 (7.74)	2.0 (0.21)	1029 (2.27)			01/10	09/19	82	340	27	N
		EnumServicesStatus	36 (1.26)	0 (0)	167 (0.37)			01/10	09/19	82	340	27	N
		OpenService	6 (0.21)	0 (0)	6 (0.01)			02/19	08/19	39	167	20	N
Driv	DeepInstinct (2019)	GetServiceDisplayName	0 (0)	0 (0)	0 (0)			Never	Never	24	87	10	N
		GetServiceKeyName	0 (0)	0 (0)	0 (0)			Never	Never	24	87	10	N
		EnumDeviceDrivers	5 (0.17)	0 (0)	21 (0.05)			01/15	08/18	60	198	16	N
		GetDeviceDriverBaseName	4 (0.14)	0 (0)	7 (0.02)			01/15	10/16	60	198	16	N

9

niques are viable only on single-core machines—they are not useful with modern systems.

Traps. Another way to throw exception is the use of x86 trap instructions. As a matter of fact, INT 3 is the typical breakpoint instruction intercepted by a debugger. If this instruction is present in a non-debugged program, the program receives the corresponding signal. A similar concept applies to VM trap instructions like VPCEXT.

Timing. These techniques contain all the methods to precisely measure time-flow and CPU clock ticks. These measurements are used by malware to spot the presence of slowdowns during the execution of malicious code. Most analysis systems have a significant impact on performance. These slowdowns are easily detectable by a malware sample. For example, it is possible to measure the number of ticks elapsed during the execution of a set of instructions using the RDTSC instruction. If this instruction triggers a virtual machine exit, the number of ticks is significantly higher.

Stalling. This category contains all the techniques that put the executable in sleep mode to avoid detection. These techniques work because sandboxes have limited time allocated for the analysis, so they could miss malicious activities of malware samples that activate after a long period. Sleeping behaviors can fool even debuggers; a loop that continuously invokes sleeping routines is frustrating to step through during manual debugging.

Human Interaction. In an automated system where thousands of malware samples are analyzed inside virtual machines, there is no human interaction. There is a category of evasion techniques that exploit this absence of interaction to identify the instrumented environment. For example, a malware sample can assume that the environment is instrumented if the mouse cursor is not moving.

Registry. Some malware samples access the Windows Registry searching for the virtual machine or debuggers-related artifacts such as *VirtualBox Guest Additions*. This class of techniques considers the use of a blacklist of resources that should not be accessed.

System Environment. This class contains all the techniques that detect system-related information that can reveal the presence of a VM or debugger. For instance, an using `GetAdaptersInfo` to retrieve the MAC information. MAC address or graphic devices have a manufacturer-fixed part that is well known for virtual machines.

WMI. An extension of the System Environment category is the WMI class. WMI is a Microsoft technology for Windows machines maintenance. Using the WMI framework and the corresponding query language WQL, it is possible to query any information about Windows machines' configuration and change their settings. Evasive malware can exploit these capabilities to detect analysis tools or installed software.

Process Environment. This category refers to techniques used to collect information about the running process. An example is `NtGetContextThread(CONTEXT_DEBUG_REGISTERS)`: this API call retrieves the values of hardware debug registers, employed to perform hardware debugging.

File System. This category includes all those techniques that rely on artifacts present on the file system (e.g., the presence of a Python-based agent, `agent.py`, in the Home folder).

List Processes. Malware can enumerate the running processes to identify an agent process (e.g., `python agent.py`).

List Services. These techniques enumerate the active services to identify an analysis service.

Drivers. Malicious samples can enumerate the list of drivers to identify emulated or virtualized devices.

Complexity of the Evasive Techniques. We study the complexity of the techniques mentioned above and the difficulty of their implementation. To do so, for each of the 92 techniques, we implement a program that leverages a given technique to detect when

it is running in a controlled environment, i.e., returning True or False. By statically analyzing each produced binary program, we count the number of basic blocks, instructions, and function calls required to implement the 92 evasive techniques (columns **BB**, **I**, and **C** in [Tables 3](#) and [4](#)). We also report whether the implementation of each technique, at the source-code level, requires writing assembly code, which indicates more advanced skills (column **ASM** in [Tables 3](#) and [4](#)). While none of the studied techniques require specific privileges or can be considered complex or challenging to implement, some groups of techniques certainly rely in smaller portions of code. For instance, Memory Fingerprinting techniques, although they require writing assembly code, are implemented with a few instructions (8 to 16) and without invoking any additional functions. On the contrary, List Processes and List Services techniques result in larger portions of code (up to 340 instructions) containing several function calls (up to 27).

4. Analysis framework

To perform our study, we developed an automated analysis framework to analyze the evasive behaviors of Windows programs. Following a recent research trend ([D'Elia et al., 2019; 2020; Polino et al., 2017](#)), which demonstrated the effectiveness of DBI tools in defeating evasion attempts, we leverage Intel Pin to fully control the execution flow of the instrumented program. In fact, by instrumenting a binary, we can identify when it attempts to spot artifacts to evade the analysis. Our framework extends a previously proposed system ([Polino et al., 2017](#)), and it can instrument executables at four levels: (1) Instruction (monitor each executed instruction); (2) API Hook (hook interesting Windows APIs); (3) Syscall (hook interesting syscalls); (4) Memory Access (monitor access to specific memory regions).

Leveraging these four instrumentation levels, we implemented a set of countermeasures that dismantle all the 92 presented evasion techniques. Our system inserts two hooks before and after instructions or function calls that can be associated with an evasive technique. The first hook is executed *before* the evasion while the second *after* it. Both hooks contain a logging routine that logs the evasive behavior and applies hiding countermeasures. We are able to insert hooks with such controls at the instruction level, leveraging Just-In-Time recompilation of Intel's PinTool. Our system has four main modules: the **System Calls Hooker** and the **Windows API Hooker**, which detect and apply the countermeasures of evasive techniques based on system calls and Windows APIs, and the **Memory Accesses Controller** and the **x86 Instructions Monitor** for evasion based on memory reads and x86 instructions. Moreover, we implemented a **Library Tracer** module to log the call stack of the evasive technique and related libraries. Windows libraries execute some instructions, API calls, or syscalls for legitimate purposes. Therefore, if we consider only their presence, we cannot claim that the malware is evasive. Using the Library tracer, we can detect which activities are executed from the malware text segment.

It is important to note that our system does not stop its analysis after the first detected technique. In the case of malware samples that implement a sequence of multiple evasions attempts, our system detects the first technique executed and applies the corresponding bypass mechanism described before, avoiding that the malware triggers evasive behavior. In this way, our system can analyze and identify any secondary techniques executed after the first "barrier."

A common bypass mechanism employed throughout our analysis consists of hiding artifacts that are commonly present in sandboxes. As an example, whenever a malware samples accesses file information executing syscalls like `NtOpenFile`, `NtCreateFile`, `NtOpenKey`, or `NtQueryAttributesFile`,

Table 5
Blacklist of functions accessing OS resources

Anti-VM
vbox, virtualbox, vmware, vmx, vmsvc_, prl, xen, parallels, bochs, wine, qemu, vmusvc_, innotek, vmhgfs
Anti-Debugging
ollydbg, debugger, debug, softice, x32dbg windbg, immunity debugger, dbgviewclass
Anti-Analysis
wireshark, tcpview, autoruns, regmon, procepx, hookexplorer, sysinspector, petools, dumpccap, python, malware, virus, sample, anyprotect, process viewer, bitdefender, seguridad, charles, dll injector, fiddler, file monitor, filemon obsidiangui, processhacker, procmon, progman regmon, rock debugger, rootkitrevealer, socketsniff, pe explorer, process explorer, process hacker, process heap viewer, process monitor, program manager, packet analyzer, registry editor, resource hacker, security task manager, switchsniffer, tcpview, wireshark, windows file protection, winhex, taskmgr, trw2000, remote process viewer, httpanalyzer, mcafee, tidawindow, smartsniff, kaspersky, norton, panda, webroot, sophos, avira, trendmicro, comodo, avira, bitdefender, clamav, symantec, avast, nod32, drweb fortinet, g data, systracer, regshot, windows task manager

our framework returns a NOT_FOUND response if the file path contains suspicious keywords defined in Table 5 (e.g., “virtualbox”). In other words, our system logs the attempt to retrieve information about the environment and ensures that the API call returned value corresponds to a not existing key. A similar countermeasure is applied when the malware accesses Keys inside the Windows Registry or enumerates drivers, Windows services, and processes. Differently, a relevant bypass mechanism involves reducing the value returned by time measurement instructions such as RDTSC to eliminate the overhead introduced by analysis environments. We include an extensive discussion of our evasion countermeasures in Appendix A. Moreover, our countermeasures follow the directions established in previous work (D’Elia et al., 2019; 2020; Polino et al., 2017).

Validation. We manually validated our framework by developing “test” programs and verifying that our framework could identify the usage of all the 92 (known) techniques described in Section 3. Moreover, we tested our detection system against (Pafish, 2020), a known tool that implements several evasive techniques. Our framework was able to detect all the techniques implemented by Pafish, remaining undetected. Naturally, this does not guarantee that our framework can detect all the malware evasive behavior that malware can adopt, as it is an undecidable problem. We discuss our limitations in Section 6.

5. Evasive behavior analysis

In this section, we describe our experimental setup and the results of our study, which specifically aims to answer the following research questions.

RQ1. *What are the most common evasive techniques adopted by malware?* We analyze 45,375 samples from 2,867 different families and report insights about the adopted evasive technique together with the specific keywords searched by samples to spot analysis artifacts, e.g., “vbox” (Section 5.2).

RQ2. *Is there a correlation between malware families and the adopted evasive techniques?* We study this correlation by training and evaluating a set of classifiers that aim at distinguishing malware families by looking at the set of adopted techniques (Section 5.3).

RQ3. *Has the number of adopted evasive techniques, per family and per category, changed over the past 10 years? When did each*

evasive technique appear for the first and last time? We perform a longitudinal study over the past 10 years, and we provide insights into the evolution of evasive behaviors over time (Section 5.4).

RQ4. *How does the adoption of malware evasion techniques impact on the security community, and vice versa?* We harvest papers, reports, and blog posts from the security community estimating the influence of the public disclosure of evasive techniques on their adoption in the wild (Section 5.5).

RQ5. *Are evasive techniques also adopted by legitimate software?* We analyze 516 legitimate programs and study evasive behaviors in goodware. We further compare such results with the data obtained from the analysis of malware samples (Section 5.6).

5.1. Datasets & setup

We leverage two datasets: a malware dataset and a goodware dataset. For the former, we collected 45,375 malware samples from VirusTotal, distributed in a timeframe of 10 years. Samples are ordered by the first appearance date on VirusTotal. In particular, we downloaded and analyzed more than 1,000 samples per quarter, from January 2010 to September 2019, in order to perform a longitudinal, time-based analysis of the evasion phenomenon. We collected only samples that were classified as malicious by at least 35 different AV scanners. Such a conservative choice is coherent with what a recent study suggested to mitigate the known label dynamics issue (Zhu et al., 2020). A lower number of AV detections would have caused the download of possible false positives, which would have biased our analysis. Nonetheless, we acknowledge that, as a consequence of this choice, our dataset might not include samples that are not well detected by the majority of endpoint solutions (Section 6).

We leveraged AVClass (Sebastián et al., 2016) to obtain the malware family of each sample. Our dataset contains samples belonging to 2,867 unique families, while 1,369 samples (3%) were not tagged in any malware family by AVClass—we included these samples in our experiments, but we did not consider them when studying per-family trends. Figure 1 shows the number of families detected in each quarter. The top 1,000 families account for almost the whole (93%) dataset, while the top 100 families account for around 67% of the dataset.

For the benign dataset, we downloaded 327 executables from PortableApps.com and 622 generic PE files from a variety of Windows OS executables. In total, we collected 949 benign executables. All the benign applications were executable desktop applications. None of our benign samples was installed from Windows Store. We filtered out any executables that were detected positives by at least one anti-malware on VirusTotal. After this process, our dataset was composed of 516 goodware binaries. All our executables were compatible with Windows 7 and later versions.

Experimental Setup. We deployed our analysis framework on 6 Windows 10 VirtualBox VMs. Each VM is managed by our analysis manager, which spawns the VM, sends the sample to the VM agent in charge of the DBI analysis, and collects the evasion report. After each execution, we rolled back the virtual machine to a clean snapshot. As done in the previous work (Polino et al., 2017), we let each sample run for up to 5 min, a reasonable time to trigger most of the evasive behaviors that our tool can detect. In fact, a recent study (Küchler et al., 2021) showed that most of the behaviors manifested by malicious samples in a sandbox (and 98% of the executed basic blocks) are observed during the first two minutes of execution. Finally, following the best practices for malware experiments (Rossow et al., 2012), we allowed the samples to communicate with their control servers and denied any potentially harmful traffic (e.g., spam) during the experiments. Our analysis machine runs on a separate local network with no other machine connected, thus, self-spreading malware does not have any local

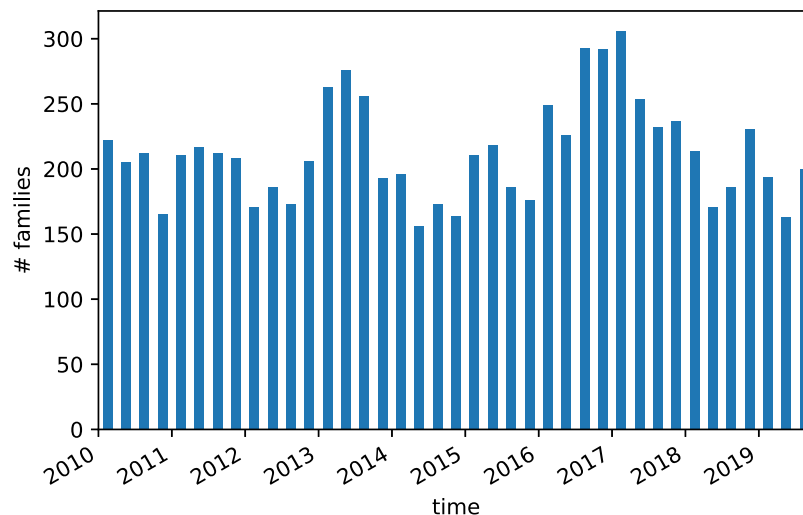


Fig. 1. Number of families over time.

target. To support scientific repeatability, we will release our analysis framework.

5.2. Adoption of evasive techniques

Among the 92 collected techniques, we measured the percentage of techniques adopted during each considered year. Tables 3–4 show statistics about the presence of each technique in our malware datasets, considering also the taxonomy introduced in Section 3. Our data show a significant usage of Timing and Stalling techniques in both malware and goodware, which use them for synchronization and timing purposes. Also, an interesting result is that 15 techniques described in the literature have never been observed in our dataset, neither by malware nor goodware samples. These techniques can be used for evasion purposes, but they do not seem actively exploited by our collected samples.

Table 2 contains data about the evasiveness of the top 20 most prevalent families. Collected metrics show that most of them first appeared in 2010 and that not always the most prominent family is the most evasive. In terms of evasion, the most evasive families are *installerex* and *domaiq*, with a striking 99% of evasive samples, immediately followed by *vtflooder* and *neshta*, with 98% of evasive samples. This aspect often pairs with the maximum number of different techniques detected in a sample of the related family. Indeed, the *installerex*, the *vtflooder*, and the *neshta* families contain samples employing at least 15 different evasive techniques. The only exception is represented by *domaiq* that, instead, employs up to 7 different techniques. Interestingly, *neshta* is also the family with the highest number of samples that are packed according to PEiD. Table 2 also shows rarely evasive families such as *xorist* and *bladabind*, with 11% and 21% of evasive samples, respectively. Another interesting takeaway is that families like *fareit* and *flystudio* have few samples exposing many evasive techniques. Figure 2 shows the normalized quarterly distribution of samples employing evasive techniques for each category. On average, the most common exploited categories are Timing and Stalling. Search for specific drivers is very rare, and we never observed the use of WMI.

Keywords Blacklist. Some techniques make access to system resources such as files, directories, Windows registries, processes, and services. Malware samples actively search for analysis and monitoring tools to evade detection. In fact, the presence of such tools is a symptom of being on an analyst's machine. Besides key-

words related to virtual machines and debuggers, we inserted keywords relative to these analysis tools (Anti-Analysis). Table 5 shows the blacklist defined for the analysis. The Anti-VM group contains some keywords related to common Virtual Environments used for the analysis. Tools considered are: Xen Hypervisors, VirtualBox, VMware, Parallels, Wine, QEMU, Cuckoo, and Bochs. The blacklist has been built selecting keywords related to VMs, debuggers, and malware analysis tools. Some of them have been collected from threat reports (Cyberbit, 2007; Project, 2017b), while we found others inside our dataset. These resources should not be visible to evasive malware because it could use them to identify our system. Indeed, we deny access to all sensible resources.

Packing. Commercial packers often employ anti-debugging techniques. To explore such a phenomenon and validate our evasion detection techniques, we ran PEiD (2020) on the samples in our dataset. PEiD detected packers in 4,293 samples, with the most common packers being UPX (1713), Nullsoft PiMP (1609), PEXcompact (237), ASPack (185). Furthermore, PEiD detected 5 samples packed with (different versions of) Themida, which is known to adopt anti-debugging techniques. For all the Themida samples, our framework detected Timing, Exception Handling, Traps, and Registry techniques. Interestingly, for versions < 1.8 we identified 3 techniques: *GetLocalTime* (Timing), *NtClose(INVALID_HANDLE)* (Exception Handling), *ANTI-VM VPCEXT* (Traps). In later versions, we additionally found *RDTSC/D* (Timing) and *NtOpenKey/Ex* (Registry).

It is important to note, the amount of packers that we identified is a *lower bound*. Indeed, PEiD is able to identify just common packers, but it does not run an entropy analysis. This implies that custom and not very widespread packers were not identified.

Techniques Clusters. We analyzed the use of techniques in malware samples by computing the Pearson correlation coefficient between each technique. We then created clusters of techniques whose correlation was above 0.5 (highly correlated). Table 6 shows the cluster of techniques that are often employed in the same malware samples.

RQ1. *The vast majority of malicious samples adopt at least one evasive technique, with Timing and Stalling being the most adopted categories of techniques. Malware families also look for artifact keywords, with “vbox”, “sandbox”, and “virtualbox” being the most used keywords.*

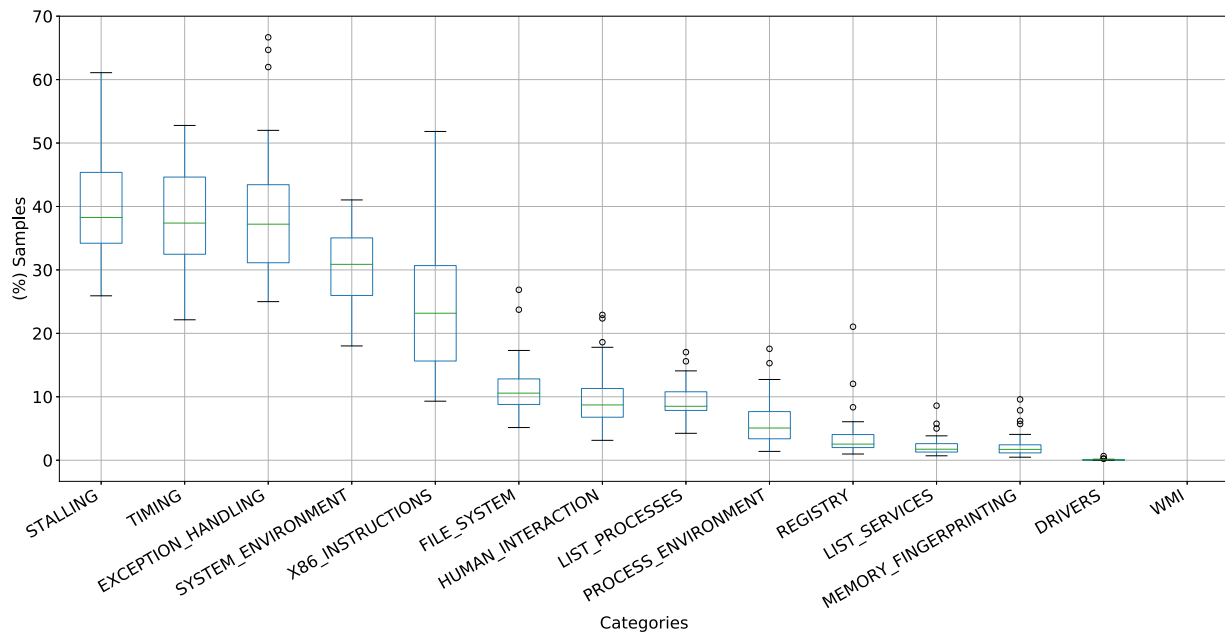


Fig. 2. Percentage of evasive samples by category.

Table 6

Techniques that are employed in the same malware samples.

```
NtQuerySystemInformation( PHYSICAL_MEMORY_INFO)
GetTickCount, RDTSC/D, CPUID(eax=0x00000001),
waitForSingleObject/Ex QueryPerformanceCounter,
SetUnhandledExceptionFilter, Sleep/SleepEx
SGDT, SLDT, SIDT
IN, VPCEXT
NtQueryInformationProcess(0x1e - 0x1f - 0x7), POP SS
CheckRemoteDebuggerPresent, NtSetInformationThread(0x11)
GetModuleName, EnumProcesses
STR, POPFD
```

5.3. Families & evasive techniques

Intuitively, the set of evasive techniques used and the malware families could be correlated. This hypothesis is raised from the way malware samples are developed. In fact, malware samples are usually built by tools. The production chain of malware has been proved to be well distinct. Malware authors usually develop builders that can be easily used to generate new samples. The malware authors benefit from the selling of the build because the criminals who distribute malware are not necessarily the ones who wrote it or even able to develop a sophisticated malware program. Evasive techniques are usually an additional feature of malware sample builders. We leveraged a machine learning algorithm of pattern identification to evaluate this correlation between the evasive techniques and malware families. We set up a RandomForest Classifier; the input of the classifier is an array of boolean variables representing the use of each technique inside a sample; the output of the classifier is the family label. We built 123 classifiers for 123 families (all families with at least 50 samples). The output of each classifier is trained with 3-fold cross-validation using stratified sampling. For each classifier, we computed the f1-score, a metric to evaluate how well the classifier is performing, taking all the confusion matrix into account. In Figure 3, we can appreciate the result of this experiment. While clearly, there are families with a solid correlation, the values of the F1-score drop quickly. This implies that from our analysis, 16 families can be recognized from the combination of evasive techniques employed. Indeed, F1-score

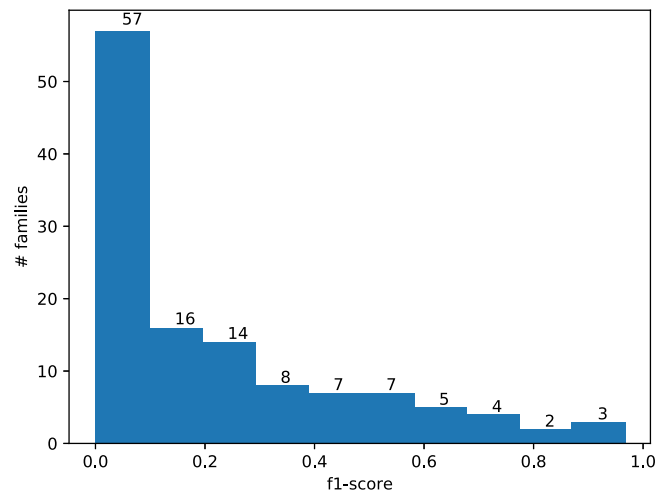


Fig. 3. F1-score for RandomForestClassifier; y-axis shows the number of correctly classified family, x-axis is ordered by f1-score.

is not enough to evaluate the performance of a classifier. An F1-score greater than 0.5 can still be produced by a low precision and recall. Figure 4 show more measurement about the performance of such classifiers. From Figure 4 we can see that two of those classifiers (imali, relevantknowledge) have low precision and recall.

Malware samples are often packed. Some packers, like obsidium, can implement evasion techniques themselves. Thus, we ran an additional experiment to identify if there is any correlation between the packer that we identified in our dataset and the techniques. The experiment setup is similar to the previous one. In this case, instead of building a classifier for the family, we built a classifier for packers. We lowered the minimum threshold of samples to 20 to include more packers. In total, we classified 33 packers; the results are shown in Fig. 5. Only 'PEtite' showed a significant correlation between used techniques and the packer (F1: 0.87, Recall: 1.0, Precision: 0.82).

Finally, while we considered tf-idf as techniques to prove correlation among techniques and families, we chose Random Forest to

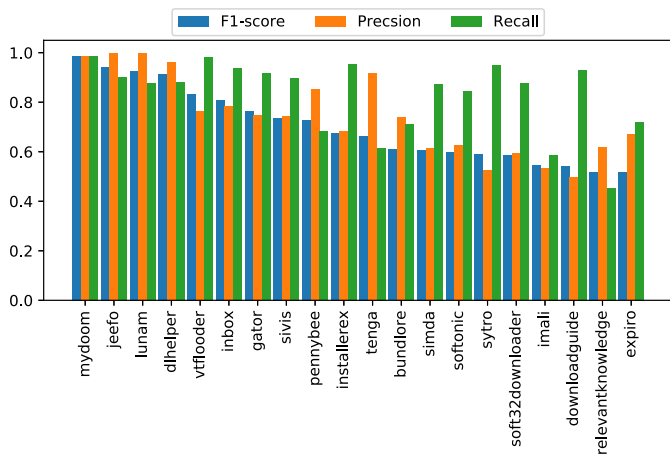


Fig. 4. x-axis shows the name of the family that was used to build the classifier. y-axis show the scores for precision, recall and f1.

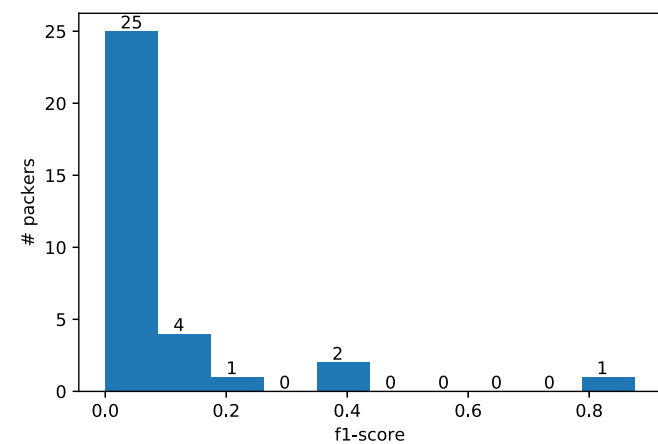


Fig. 5. F1-score for RandomForestClassifier; y-axis shows the number of correctly classified packers, x-axis is ordered by F-1 score. The only packer with F-1 score greater than 0.5 is PEtite.

study further if a detector with good performance could be built for some of the families and packers.

RQ2. There is a consistent number of malware families (46) that are characterized by the selection of evasive techniques they employ.

5.4. Longitudinal analysis

Our dataset has been collected evenly spread across ten years (Section 5.1). This dataset empowered us to perform analysis and identify changing trends during these years.

Interestingly, Fig. 6 shows that there was a little increment (~ 12%) in the use of evasive techniques during the last ten years. In comparison, the use of common packers stayed constant. The detection of all techniques starts from 67.8% in 2010 and slightly increases to 80% during the last years. This result is coherent with old works (Branco et al., 2012; Chen et al., 2016). Although, in Section 5.6, we discussed how several techniques could be confused with normal program behavior.

Moreover, we analyzed the intricacy of malware during the last ten years. Figure 7 shows the maximum number of evasion techniques that were used by a single sample. It is interesting to notice the rising of the number of techniques until 2017. However, we can see a decrease in the last two years.

Furthermore, we analyzed the targets of evasive attempts. Indeed, our system can discern when the malware is evading a de-

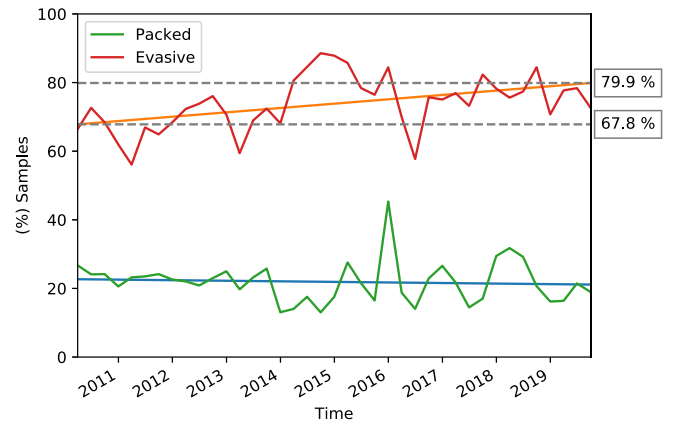


Fig. 6. Percentage of evasive samples per quarter. Straight lines are linear regressions: In orange for evasive techniques, in blue for packers. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

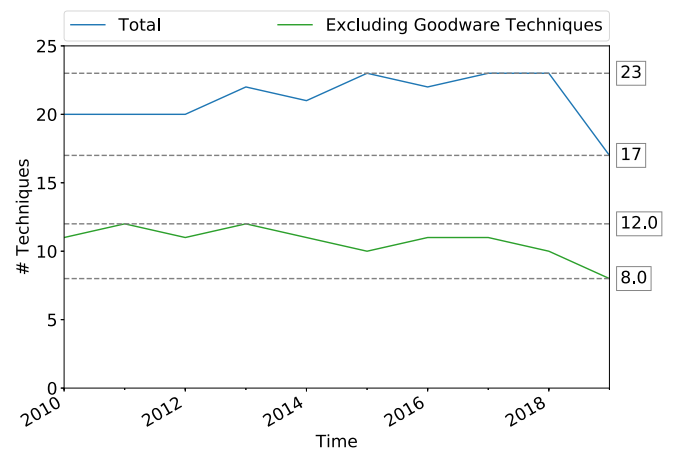


Fig. 7. The maximum number of evasive techniques detected in a single sample, yearly aggregated.

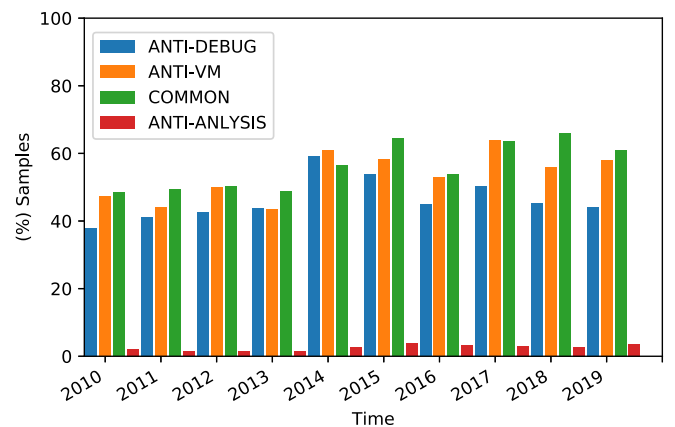


Fig. 8. Percentage of samples divided by targets.

bugger or a virtual machine. In Fig. 8 we show four categories of targets.

We exploited the categorization in Tables 3-4. However, for techniques used for both Virtual Machines and Debuggers, we relied on the blacklist described in Section 5.2. There are few techniques categories (Stalling, Timing, Process Enumeration) that cannot be assigned to either of the categories. We named these tech-

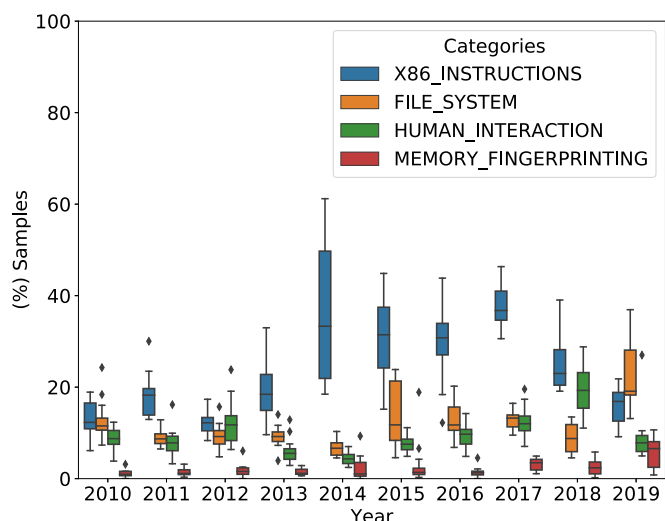


Fig. 9. Timeline of techniques usage.

Table 7
Usage of blacklisted keyword per families in our dataset

Family	keyword	#
neshta	sandbox	444
	virtualbox	443
	xen, sample	442
	python	424
	debug	353
	vbox	226
	malware	203
outbrowse	vbox	302
installmonster	virtualbox	129
sivis	debugger	127
hotbar	debug	118
simda	wireshark, virtualbox, debug, vbox	85

niques *Common*. Figure 9 shows how the type of techniques used by malware authors shifted during the last 10 years.

Our blacklist includes, besides virtual machine and debugger, generic analysis tools keywords. In fact, finding classical analysis tools is a clue that the machine is used for analysis. The *Anti-Analysis* peak in 2018 is caused by a campaign of the *neshta* family that scans the disk searching for files to infect. In 2015 instead, the predominant *Anti-Analysis* families were *simbda* looking for *wireshark* and *eorezo* searching for antivirus software. Keywords used for most common families are shown in Table 7.

Another interesting aspect is the timeline of the techniques: lots of them have been observed during the whole timeframe considered, which means that their first appearance may go back to years before the ones considered in this work and that malware authors will probably continue to use them in the future. Figure 10 shows the period of utilization of each technique. We restricted the graph to techniques that are predominantly found in malware (See Section 5.6). The graph is based on the observation of running samples, and the dates refer to the first appearance of the corresponding sample on VirusTotal. Several techniques start at the beginning of 2010. Very likely, those techniques have been used before 2010. The same reasoning can apply to the ending: the closer the ending date is to the end of our dataset (September 2019) more likely it is that the technique is going to be used in the future. We plotted each observation – of an evasive technique – as a vertical line on the corresponding bar. We can notice two patterns: (1) many techniques shows cluster in small periods. This can be attributed to active malware campaigns, either because analysts submit newly found samples or because malware authors test their samples before release. (2) Some tech-

niques are more used than others. We can notice like both *INT 2D* and *IsDebuggerPresent* are really old, but the latter is definitely more used. Indeed, *IsDebuggerPresent* is documented by Microsoft API documentation and technically easier to implement; these may have contributed to its popularity as an evasion technique.

Furthermore, we can notice techniques that are not used anymore, like *PEB->NumProcessors* that checks the number of cores and *SIDT* that exploits a side-effect of virtualization techniques in old CPUs. Those techniques are not effective anymore since the rise of the number of cores in modern processors allowed virtual machines to use more than a single core and modern CPUs support for virtualization removed some of the side-effects.

RQ3. Our longitudinal study highlights that evasive techniques were already quite common back in 2010, with almost 70% of malware samples employing evasive behaviors, and their adoption slightly increased over the past ten years, although not significantly. Moreover, we can see malware authors drop techniques that are not effective anymore.

5.5. Evasion vs. community

Evasive techniques have been reported for years in research papers, industrial conferences, and specialized forums. In this section, we conduct a preliminary study on the impact of the adoption of malware evasion techniques on the security community and vice versa. In particular, we aim at estimating (1) the influence of the release of information about evasive malware techniques on their adoption and (2) the time required by the community to report a new evasive technique with respect to its first appearance in our dataset.

We crawled the Internet for information related to the evasive techniques under analysis to estimate the community awareness on malware evasiveness. We focused on academic and industrial works scraping paper, blog posts, and technical reports. We searched for the keywords related to the techniques listed in Tables 3 and 4. Then, we filtered the results associated with the concept of malware or evasion (e.g., “*IsDebuggerPresent*” AND (malware OR evasion)).¹ We further post-processed the gathered data by manually analyzing and filtering out non-relevant results (i.e., generic reports of dynamic analyses with no explicit references to evasive behaviors). Coherently with the longitudinal study conducted in this paper, we focused on the period from 2010 to 2019. Additionally, we crawled the same information for the period before 2010 to study the impact of the security community on already known techniques. We collected a total of 394 sources². Then, we compare each evasion technique’s quarterly trends (i.e., the number of sample implementing it) with the harvested information. Despite this study provides only an estimation of the influence of the community on the adoption of evasive techniques – and vice versa – since it might be biased by both the dataset under analysis and the (publicly available) sources, we deem that it provides insights for future research on the topic. In fact, in the period under analysis, we observe an increasing trend in the number of resources related to malware evasiveness: between 2010 and 2015, we have an average of 3 reports per year, while between 2016 and 2019, this number doubled, demonstrating the higher community interest and awareness of the problem, passing the years. From the 92 evasive techniques under analysis, 50 were already known before 2010, while 21 are disclosed by the community between 2010 and 2020. On average, techniques become public after 4.4 years. The fastest disclosure belongs to the ‘*EnumDeviceDrivers*’

¹ To limit the bias of search engines’ personalized search and geolocalization, we worked in “Incognito mode”.

² The gathered data is available at <https://github.com/necst/brioscia>

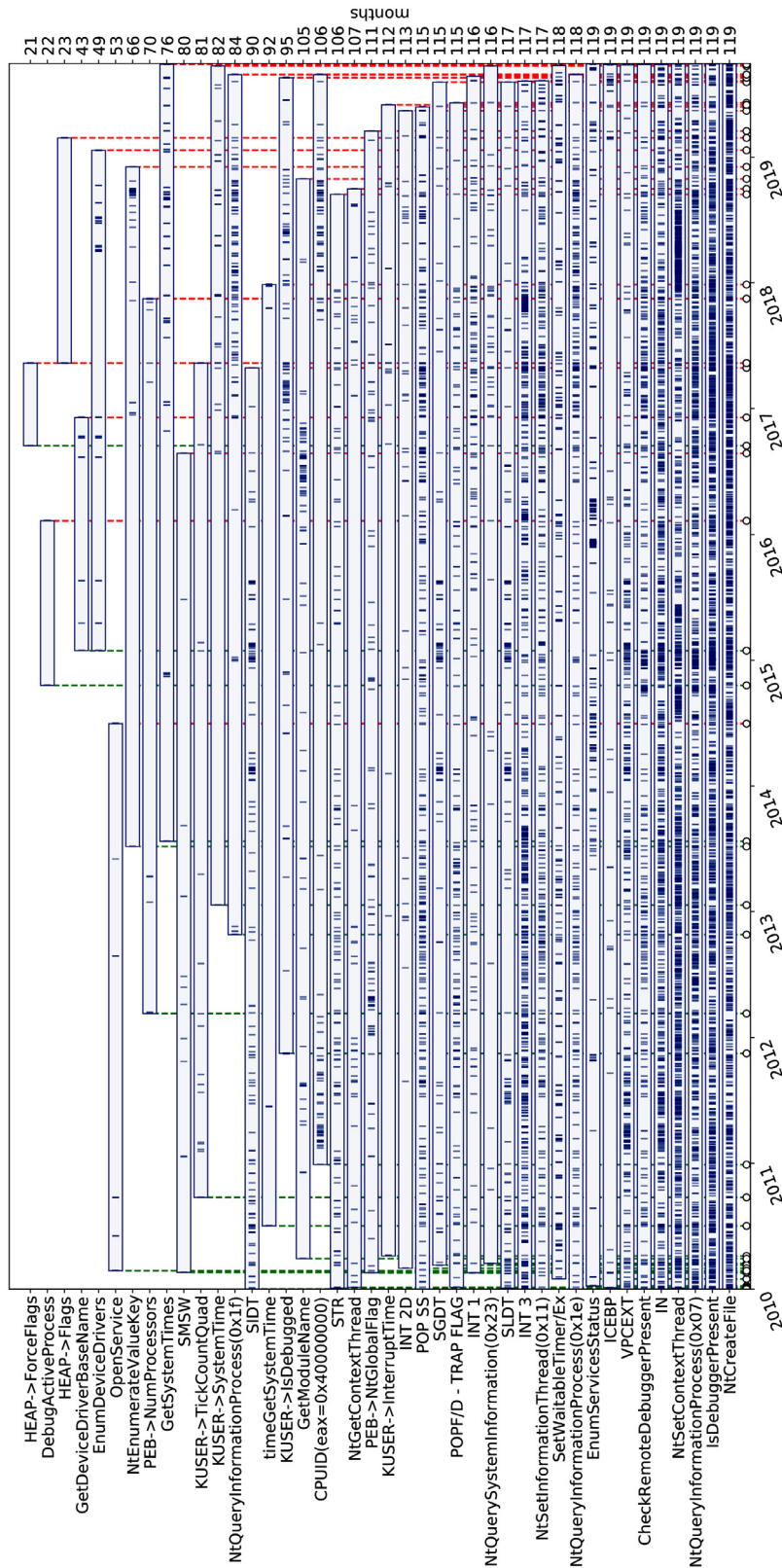


Fig. 10. Time-line of techniques usage in malware samples. Each horizontal line indicates the period in which the technique has been observed. Each tick is one observation of a sample using the technique.

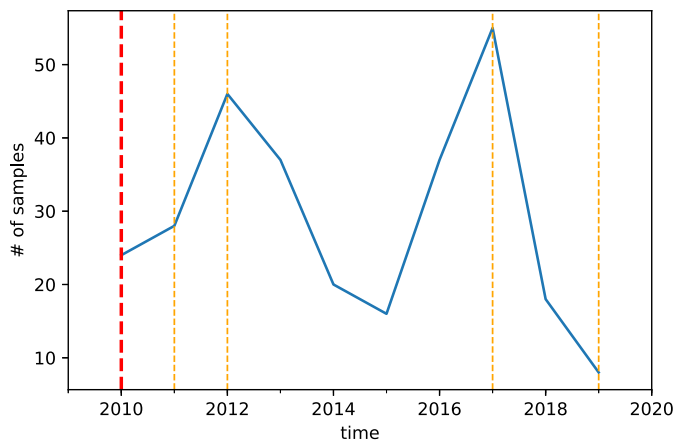


Fig. 11. POP SS detection timeline. The solid line is the number of samples that adopt POP SS. The thick dashed vertical line is the first public appearance of the technique, the thin dashed lines any sequent appearance.

technique, which appears in malware samples in the same quarter of the talk (Timzen, 2020). Instead, the slowest disclosure—after 7 years—belongs to ‘VPCEXT’ (McAfee, 2017).

By comparing the quarterly trends of each evasion technique with the harvested information, as shown in Fig. 11 for the POP SS technique, it is possible to observe an evident influence between them. There is a high correlation between the release of information about evasive techniques and their implementation in malware. In particular, if the technique is already widely adopted/known, after the release of new information, we notice a slight and temporary decrease in their usage in the subsequent years. On the contrary, if the technique is scarcely used/known, we observe a swift increase in its adoption in the following years after the release of reports. As expected, once a technique becomes widely used, the number of related reports significantly increases.

RQ4. *Our analysis highlights a correlation between the adoption of malware evasive techniques and the information provided by the security community: if from one side, we observed that the release of security reports on evasive techniques corresponds to an increase of their adoption in the short term, on the other, it tunes down their spread on the long term.*

5.6. Malware vs. goodwill

While less common, goodwill also employs evasion techniques to prevent reversing and protect intellectual property (e.g., games, Spotify). Moreover, some of the evasive techniques that we monitor may be based on side-effects and operations that are not intrinsically malicious, and that can be used both for evasive and non-evasive behaviors. For instance, Operamail monitors the list of running processes to look for another instance of itself, and it does not run if another instance is already running.

In order to better characterize such techniques, we collected a dataset of benign applications. We discarded all those binaries that were positively detected by any anti-malware used by VirusTotal. The goal is to distinguish techniques that are used exclusively for evasion purposes from the ones that are, instead, often used by legitimate programs for benign behaviors. For instance, CPUID(`eax=0x1`) is an instruction that Microsoft compiler adds at the beginning of any binary. While this instruction can be used to fingerprint the CPU and be used to understand if the program is executed inside a virtual machine, it can also be used to understand if a library needs to use a software implementation of floating-point arithmetics or it can exploit the hardware FPU.

Table 8

Techniques that are employed in the same goodwill sample.

```

GetCursorPosition, RDTSC/D,
GetTickCount, SetUnhandledExceptionFilter,
NtQuerySystemInformation(PHYSICAL_MEMORY_INFO)
NtQuerySystemInformation(SYSTEM_PROCESS_INFORMATION),
GetAdaptersAddresses, EnumProcesses, timeGetTime
IsDebuggerPresent, NtQueryValueKey,
NtOpenKey/Ex, GetDiskFreeSpace
NtCreateFile, NtOpenFile
NtDelayExecution, Sleep/SleepEx
NtQuerySystemTime, GetComputerName,
NtClose(INVALID_HANDLE)
OutputDebugString, PEB→IsDebugged

```

We then analyzed these applications to identify evasive patterns. Figure 12 shows the number of evasive techniques adopted by legitimate programs in our goodwill dataset (30). We can see that some evasive patterns are present in many goodwill programs. Sometimes, goodwill implement anti-analysis techniques to protect Intellectual Property. For this reason, we manually analyzed the use of 18 techniques. We choose to analyze the techniques that are less frequently used in benign programs. One (`IsDebuggerPresent`) of these eighteen techniques were used as an evasion technique to identify the presence of a debugger. The remaining 17 were not used as an evasion mechanism. We classified each technique based on our goodwill analysis:

Used by Malware. It is the evasive behavior never observed inside goodwill samples, or we manually verified that it is used as an evasion mechanism in goodwill samples.

Used by Goodwill. It is the evasive behavior observed in more than 1% of our goodwill samples, or we manually verified that there were not used as an evasive behavior.

If we consider this division of techniques and revisit some of the results we have gotten so far, we can notice some differences. Indeed, looking at Fig. 13, we can see that if we consider only malware-related techniques, the percentages of evasive families drop from 60-85% to 19-30%. While the use of techniques used by goodwill is increasing, the use of techniques that are explicitly only in malware stays constant. Similarly, looking at Fig. 7, if we consider all techniques, there is an increase in usage until 2017. If we take into account only techniques found exclusively in malware, the trend shows a constant decrease in the number of techniques used in the same sample. We also remark that techniques never detected in benign programs are the ones used explicitly for malicious purposes. Even though we validate our benign dataset against VirusTotal to minimize the presence of trojan in our benign dataset, it is possible that some Trojan was not caught by our validation. The lower bound results may be influenced by those trojans.

Techniques Clusters Similarly, for what we did for malware techniques, we computed the Pearson correlation among techniques that are implemented in the same goodwill. Techniques with correlation > 0.5 are shown in Table 8. Interestingly the first group shares all but one (`GetCursorPosition`) technique with the first group of malware samples (Table 6). Besides that, all other clusters are different from the ones seen for malware samples.

RQ5. *Legitimate programs adopt evasive techniques, even if this phenomenon is significantly less common. In particular, only a few techniques are commonly adopted, suggesting that such techniques are not mainly used for true evasive purposes.*

6. Limitations

Reliable detection of evasive techniques is a complex engineering problem, as a “perfect” solution often does not exist, and the adversarial context exacerbates the issue. We acknowledge that

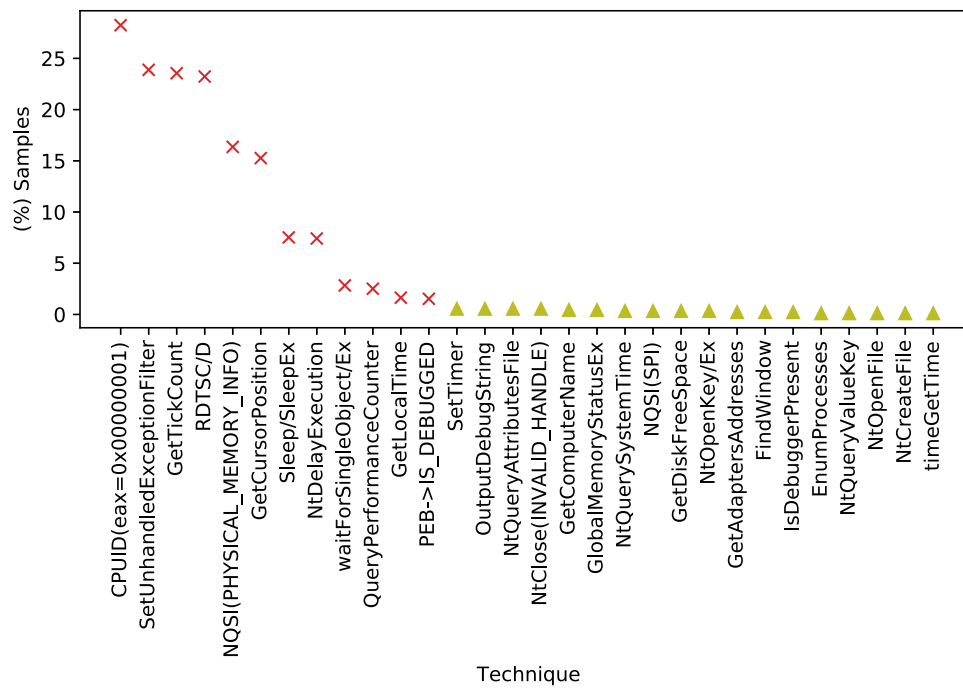


Fig. 12. Appearance of evasive techniques in Goodware dataset. × markers identify unreliable techniques because of their massive presence in goodware programs. △ markers identify techniques that we have manually analyzed in goodware programs to determine if the technique is used for evasion.

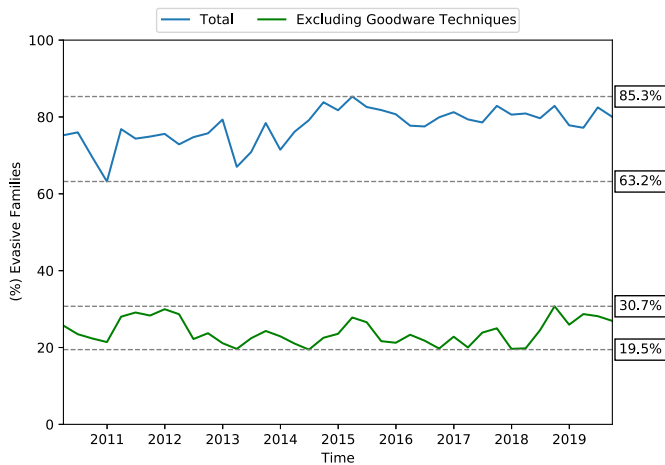


Fig. 13. Ratio of evasive families quarterly aggregated.

our analysis framework itself can be evaded: We measure the usage of the 92 (known) techniques that our framework supports, but of course, samples may employ unknown techniques that we are unable to track.

Our choice of selecting samples detected by at least 35 AV engines allows us to exclude false positives (i.e., legitimate programs that we would consider malicious). However, as a side effect, our dataset may not include true malicious samples that are not well detected by the majority of AV engines. Unfortunately, this is an inevitable trade-off since manual verification of each sample is not feasible when performing large-scale studies such as ours. We plan to investigate evasive behaviors in less detected samples in future research.

Similarly, to avoid malicious samples into the benign dataset, we filtered out any binary that was detected as positive by an AV

on VirusTotal. Unfortunately, we cannot be certain that our dataset does not contain any Trojan.

Coherently with the research aim to perform a longitudinal analysis, during data collection, we cared about temporal distribution of the number of samples trying to sample the same amount of binaries each month. Of course, this affects the distribution of samples per family. As we can see in Table 2, even in the top 20 most common families in our dataset, there is a significant difference in the number of samples. In fact, the type of families in a month depends on the currently active malware campaigns and samples submitted to VT. Moreover, this also depends on the type of family itself. In fact, if the malware sample is polymorphic or metamorphic, the number of samples is very high. Instead, for non-morphic malware, we see only one sample. This implies that when we look at results by the number of samples, the results may be affected by the most prominent families. For these reasons, we also provide, when possible, results based on the number of families.

Some technical limitations of our current implementation stem from the use of a DBI framework for the analyzer. The stability of the injection and the performance impact of the DBI tool are significant challenges. Also, the artifacts induced by the DBI tool itself could be exploited by new evasive techniques (Polino et al., 2017). This means that our results might underestimate malware evasiveness.

Moreover, any evasion attempt that is performed after 5 min is lost by our framework. If we consider the most evasive families (>95%) and we assume that those samples were misclassified by our framework, our framework lost evasion techniques for at least 121 samples. Although, accordingly with Kuchler et al. (2021), most malware samples execute within 2 min timeframe. This implies that the execution time should not have a significant effect on our results.

The techniques we monitor are indicators of potential evasion attempts, but they do not capture the intention of malware au-

thors. We can only speculate why malware authors are adopting the techniques that we identify in our study.

7. Discussion and future directions

In this paper, we shed light on the evasive behaviors employed by modern Windows malware by systematically documenting and classifying 92 evasion techniques and by measuring their adoption over the past 10 years.

From a longitudinal, historical perspective, our results show that, overall, the evasion rate had a small increase over the years, starting very high since 2010 (when it was already almost 70%) and reaching 80% at the end of 2019. We also documented a decrease in the intricacy of malicious techniques over the last years, with Timing techniques increasingly becoming their first choice. We also found an interesting relationship with publications of techniques by the research and anti-malware community, with an immediate increase in adoption, followed by a steeper decline.

While we can only speculate the reason behind the malware authors' choices, we believe that the main reason for the observed increase in the adoption of evasive behaviors is the fact that the studied techniques are still quite effective for the cybercriminals' purposes (i.e., evading dynamic analysis and delaying the characterization of their samples). In spite of the considerable amount of research proposed in the past 10 years, mitigation of evasive malware in the real world is still a problem. While, on the one hand, designing an "invisible" dynamic analysis environment is an extremely challenging task, which can be linked to the halting problem, on the other hand, we believe that the security community can and should research novel and practical solutions to mitigate this problem further. In particular, while online sandboxes nowadays implement many anti-evasion strategies, not many endpoint anti-malware solutions employ dynamic-analysis approaches (Quarta et al., 2018). Those that do adopt lightweight solutions are, unfortunately, incompatible with the resource-hungry anti-evasion approaches currently available. To overcome this, vendors often send suspicious executables to cloud-based engines for more complex analyses, but this leaves open a window of time for the execution of malicious payloads on endpoints, which in the case of ransomware may be enough to make the subsequent cloud-enabled detection pointless. Thus, evasive behaviors allow malware authors to "delay" the identification of their samples to a point when they already caused substantial damage. We believe that future research should focus on re-designing approaches that mitigate evasive malware to make such solutions practical and efficient so that they can be deployed inside emulators that run directly on the endpoints, allowing for faster detection.

Our study also found that malware authors move from old techniques to new ones, pursuing their goal to evade a dynamic analysis system. Therefore, it is essential for forward-looking anti-evasion solutions to be able to capture the fundamental properties of evasive behavior and to adapt to emerging techniques quickly.

Moreover, our analysis framework—similarly to the previously proposed anti-evasion tools—currently does not identify the adoption of logic bombs for evasion purposes (Fratantonio et al., 2016), which we plan to study in future work.

An additionally interesting future work is the analysis of evasive techniques targeting the 64-bit x86 architecture. Our analysis in this paper is based on binaries using the 32-bit x86 instruction set architecture because it is the architecture prevalently used by malware authors nowadays. However, there are a few differences in its 64-bit counterpart that may change how some of the evasive techniques work.

Finally, no previous study focused on evaluating the effectiveness of the currently deployed anti-evasion solutions. In fact, some

existing and old techniques are currently still being adopted, which suggests that they are still effective. We believe that future work should research novel metrics and principled approaches to test and compare sandboxes and instrumentation frameworks, ultimately assessing and quantifying their fingerprintability and robustness to evasiveness.

8. Conclusions

In this paper, we systematically documented 92 evasion techniques (to our knowledge, the broadest collection so far in literature) adopted by modern Windows malware and created a meaningful taxonomy of them. We implemented and released a DBI-based tool to analyze Windows executables and identify whether they employ any of such evasive techniques. Leveraging our analysis tool and a dataset of 45,375 malware samples from 2,867 different families observed in the wild over a span of 10 years, we performed a measurement of the evasion phenomenon over time. We showed that the adoption of evasive techniques slightly increased over the years and is nowadays common in modern malware and that malware authors update their techniques to evade analysis systems. From a non-temporal perspective, we found a high correlation between sets of malware families and associated evasion techniques. We also identified 15 published techniques that do not seem to be prevalent in our dataset: we will investigate this result more in-depth in future extensions of this work. Finally, we analyzed the prevalence of (some) evasive behaviors in legitimate executables by running a comparative analysis on a dataset of goodware samples. This is the first work that shows a large collection of evasive behaviors and a quantitative analysis of their distribution over time. We believe it provides interesting insights to explore other methodologies to detect novel classes of evasive behavior.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

CRediT authorship contribution statement

Nicola Galloro: Methodology, Software, Data curation, Writing – original draft, Visualization, Investigation. **Mario Polino:** Conceptualization, Methodology, Validation, Formal analysis, Investigation, Data curation, Resources, Writing – review & editing, Supervision, Visualization. **Michele Carminati:** Writing – review & editing, Supervision, Formal analysis, Validation, Visualization. **Andrea Continella:** Writing – review & editing, Visualization, Formal analysis, Validation, Resources. **Stefano Zanero:** Writing – review & editing, Project administration, Funding acquisition, Supervision.

Acknowledgements

We would like to thank our reviewers for their valuable comments and input to improve our paper. We would also like to thank VirusTotal for providing us access to some of malware samples used for this project. This work was partially supported by research funding provided by BVTech SpA.

Appendix A. Evasive techniques: details

A1. Memory fingerprinting

A debugger usually leaves some traces in the memory of the debugged process. We consider some memory locations that are interesting for Anti-Debug checks. Some of them are inside

the Process Environment Block (PEB) or the `KUSER_SHARED`³, which are Windows data structures inside every running process. Our *Memory Controller* module handles all the memory accesses. We consider the access evasive only if the load instruction is executed inside the main binary code section to filter relevant results. In particular, we consider the following memory accesses: `PEB→IsDebugged` is a byte inside the PEB, checked by `IsDebuggerPresent` (Kulchytsky and Kukoba, 2019; Xu Chen et al., 2008), which is 1 if the process is inside a local debugger and 0 otherwise; `PEB→NtGlobalFlag` is a byte inside the PEB at offset 0x68 and it is 0 when no debugger is in place and 0x70 otherwise (Kulchytsky and Kukoba, 2019); `PEB→Heap→Flags` is a byte inside the Process Heap (offset 0x18 from PEB) and located at offset 0x40 from the Process Heap, which is always different from 2, when under a debugger (Kulchytsky and Kukoba, 2019); `PEB→Heap→ForceFlags` is a byte inside the Process Heap (offset 0x18 from PEB) and located at offset 0x44 from the Process Heap, which is always different from zero when under a debugger (Kulchytsky and Kukoba, 2019); `KUSER_SHARED→KdDebuggerEnabled` is a byte inside `KUSER_SHARED`, which is 1 when a kernel debugger is detected at boot time, 0 otherwise.

A2. Error/Exception handling

Malware samples can exploit exceptions and errors to identify the presence of a debugger. Indeed, a malware sample sets a custom handler, then throws an exception, and it executes its malicious activities inside the unhandled handler. There are several ways to implement such behavior (Ferrie, 2011). Debuggers usually catch `CTRL+C` exception using their own handles (**CTRLException Handling**). Therefore, a debugger can be detected throwing a `CTRL+C` exception after having set a handler for the exception. To detect this behavior, our tool checks if a `SetControlCtrlHandler` is called before throwing this exception using `GenerateConsoleCtrlEvent`. Similarly, `SetUnhandledExceptionFilter` is used to set the last handler in the context of the Structured Exception Handling (SEH) mechanism in Windows. The *Unhandled Exception Handler* is used to perform the last recovery action before killing the process. Windows uses the standard Windows Frame as the last exception handler. However, debuggers ignore this exception to keep processes running even when an exception occurs. Moreover, a malware sample can exploit the close of an invalid handle (`NtClose(INVALID_HANDLE)`) that throws the exception with code (0xc0000008). Slightly different is the use of `OutputDebugString`. This technique does not work on modern Windows systems; however, it was used in the past. It consists of a call to `OutputDebugString`, which sets an error if a debug is not present. It is possible to recover information about the error using `GetLastError` function. In the past, `OutputDebugString` was used to crash debuggers by using a `%s` format string as an argument. Our tool detects any `OutputDebugString` call during the execution.

A3. CPU fingerprinting

`CPUID` is an x86 instruction used to retrieve information about the running CPU. It is commonly used to retrieve information about the CPU: supported extensions, CPU brand, Floating-Point unit and cache description, CPU serial number, and other information. This instruction is generally used by computation-intensive

programs that need to know, for example, if SSE and floating-point operations are supported natively. There are two common ways to detect a hypervisor using this instruction (Cyberbit, 2007): `CPUID(EAX=0x00000001)` returns information about the presence of a Hypervisor. Instead, `CPUID(EAX=0x40000000)` returns the brand of the hypervisor, such as *VMware* or *VBOX/VBOX*.

`IN` is a privileged x86 instruction used to communicate with peripherals and takes as arguments `EAX` and an immediate or `EAX` and `EDX`. Inside some hypervisors, the instruction can be executed in user-mode with some specific port values (`EDX`) to communicate with the hypervisor, which sets some special I/O ports available to guests. VMware supports this feature, indeed, `IN(EAX='VMXh', EDX='VX')` results into `EBX='VMXh'` if executed under a VMware Hypervisor, thus this method actively exploited by evasive malware to determine if they are under a VM (Cyberbit, 2007).

`SMSW` has an undocumented bits for the returned value. Some VMware versions have a default value for these undocumented bits (Sikorski and Honig, 2012).

A4. Table descriptors

These techniques find discrepancies between OS descriptor tables addresses in order to detect the hypervisor presence.

`STR` returns the pointer to the Task State Segment (TSS) of the current task. Its value is fixed on single-core physical machines but changes in virtual machines. Moreover, `SIDT` is an instruction used to retrieve the Interrupt Descriptor Table (IDT). Old VMware did not emulate this instruction correctly, so the value of the Hypervisor IDT could be leaked. In particular, if a single CPU machine is in place and a VMware Hypervisor is present, it is possible to notice two different IDT addresses for the same CPU. These methods do not work reliably on multi-CPU systems (Sikorski and Honig, 2012).

`SGDT` is used to retrieve the Global Descriptor Table (GDT), which should be unique; however, while running on VMware Hypervisor, the value keeps changing. The difference can be detected querying multiple times `SGDT` (Sikorski and Honig, 2012). Instead, `SLDT` is used to retrieve the Local Descriptor Table (LDT). This value is always zero on physical Windows machines because the Windows OS never uses it. Although, while running on VMware, `LDT` value changes (Sikorski and Honig, 2012).

A5. Traps

Several x86 instructions can be used to leak the presence of a hypervisor, emulator, or debugger (Ferrie, 2011). Similarly, to Exception Handling (see Section A.2) we can exploit `INT 3` that throws a `BREAKPOINT_EXCEPTION` that is used by debuggers to set breakpoints inside a program, or `INT 1` that is used for a `SINGLE_STEP` exception (Chen et al., 2016). `ICEBP` (opcode 0xf1) is an undocumented x86 instruction that trigger `SINGLE_STEP` exception (Sikorski and Honig, 2012). Another way to trigger exception is to set the trap flag inside `EEFLAGS` register. In fact, it is possible to pop a new value for `EEFLAGS` using `POPF` and `POPF` instructions (Kulchytsky and Kukoba, 2019). Other instructions like `POP SS`, instead, tampers with the step execution of debuggers skipping the next instruction (Kulchytsky and Kukoba, 2019). `INT 2D` throws a `BREAKPOINT_EXCEPTION`, and in some debugger (OllyDBG), this instruction will skip the one byte causing a different instruction to be executed.

`VPCEXT`, instead, is an instruction that can be executed only by Virtual PC Emulator, and it throws an exception if executed outside a Virtual PC. Some malware samples execute `VPCEXT` and then check the value of `EBX` to understand if Virtual PC detection succeeded (vpc, 2017). Virtual PC is right now an abandoned

³ `KUSER_SHARED` is a memory region that contains data shared between kernel and userspace (Microsoft, 2020).

project⁴(last version released in 2011). Thus this technique is unreliable on modern systems. For all these techniques, our system detects when corresponding instruction is executed.

A6. Timing

It is possible to exploit accurate timestamps to identify analysis systems. For example, discrepancies in the execution time of instructions are clear signs of a virtual environment. However, these techniques can also be used to detect debuggers and DBI tools (Polino et al., 2017). Indeed, any operation that slows down the execution of the sample can be detected.

RDTSC is an x86 instruction used to retrieve the number of clock cycles since reset. It is a 64-bit value used for performance measurements. In order to identify the presence of a Hypervisor, it is sufficient to execute VM Exit instructions such as CPUID while measuring the timing with RDTSC. We can detect the Virtual Machine by looking at the number of cycles needed to execute such instructions. RDTSC can also be used to measure the time elapsed between two blocks of code for debuggers or DBI detection (Sikorski and Honig, 2012).

We implemented a system to bypass our VM detection by saving the first RDTSC return value. All the subsequent calls are set to a small increment with respect to the previously saved value.

There is a series of Windows API calls that can be used for timing detection: `QueryPerformanceCounter` and `NtQueryPerformanceCounter` are unbiased versions of RDTSC that take into account multiprocessors and frequency changes across CPUs; `timeGetTime`, `timeGetSystemTime`, and `GetSystemTimes` return the time since system boot; `NtQuerySystemTime` and `GetLocalTime` return the current date in different formats; `GetTickCount`, `GetTickCount64`, and `NtGetTickCount` return the number of milliseconds since boot; `QueryInterruptTime`, `QueryInterruptTimePrecise`, `QueryUnbiasedInterruptTime`, and `QueryUnbiasedInterruptTimePrecise` return the number of 100 nanoseconds units since system boot using clock interrupts. All these functions are hooked and logged by our detection system. However, to bypass the detection of an analysis environment, all the internal Windows timers should be set properly.

`GetTickCount` and `timeGetTime` can be used to detect how much time the system has been up, and some malware sample uses them to detect the uptime of the whole system. Indeed, a small uptime is considered an indication of Sandbox (McAfee, 2019). Initially, our detection system sets timing variables to a value compatible with 7.7 h of uptime and increases their value by a small constant for any subsequent timing information retrieval.

Timing evasion techniques usually involve the repeated measurement of performance counters and ticks to check if a significant amount of time is spent on some instructions. If the delta between two measurements is less than a threshold, the system is deemed clean. Setting always a small difference between measurements bypasses all the possible timing checks.

In order to avoid the manual implementation for all the timing functions, we applied a bypass based on monitoring access to all the timing fields contained in `KUSER_SHARED:TickCountQuad` contains the value of ms elapsed since the system boot. This value is used to produce `GetTickCount`, `GetTickCount64`, and `NtGetTickCount`. The initial value is set to ≈ 28000000 ticks (≈ 7.7 h); `InterruptTime` is the field accessed by `timeGetTime`

and `QueryInterruptTime` and it is a more precise way to get the ticks of the system. This field is initialized with the value of `TickCountQuad` and incremented with a higher frequency: `SystemTime` is the field containing the value of the system date with millisecond precision. Each time any of the previous timing fields are accessed, our tool adds a tiny constant to increase the timer. Using this solution, a malware sample that continuously accesses timing information cannot reliably tell the difference between physical and virtual environments.

A7. Stalling

Malware Sandboxes are designed to detonate malware inside the Virtual Environment in order to collect system events. However, this activity is limited by the time assigned to the analysis. A reasonable time assigned by Sandboxes ranges from 3 min to 10 min, thus, a safe choice is to start malicious actions after 10 min in order to fool the detection system. This behavior is extremely stealthy. Understanding if a sample is stalling is an undecidable problem. A malware sample does not stall only using sleeping functions offered by the OS, but also performing useless arithmetic operations, or system calls (Kolbitsch et al., 2011; Oyama, 2018a). Even worse is the case of malware that contains time bombs that run only at specific dates or timestamps. Our system is limited, focusing only on the functions provided by Windows. The main sleeping function is `Sleep`, which stalls the execution of a program. `Sleep` internally calls `SleepEx` which in turn executes the syscall `NtDelayExecution`. Moreover, we can detect the use of: `waitForSingleObject/Ex`, used to wait until an event is signaled. However, it is possible to set a maximum waiting time, which causes an infinite stalling if the signal is never sent; `SetTimer`, `timeSetEvent`, `SetWaitableTimer`, used to set the execution of a function after the timeframe is elapsed; `CreateTimerQueueTimer` creates a timer inside a Queue of timers. When the timer expires a callback function is executed; `IcmpSendEcho/2/Ex` is used to wait for the response of ICMP echo responses with a timeout. If the server never replies, the function will stall until the timer expires.

We set to 0 the sleeping time defined as a parameter of the aforementioned functions. This countermeasure is applied if the function is invoked by the text segment of the malware executable because most of the sleeping functions such as `waitForSingleObject` and `SetTimer` are used for synchronization purposes by external libraries whose functionalities could be broken with a null value.

A8. Human interaction

The idea of this set of techniques is to detect if a human is using the infected machine. A Virtual Environment can be detected checking the mouse cursor position over time using `GetCursorPos` API call (**Cursor Position**). If the position of the cursor never changes, a malware sample can reliably say that it is running inside a Sandbox (FireEye, 2014). Our tool detects this API call and randomizes the coordinates to bypass the evasive behavior.

Another way to detect the presence of the user on the system is to retrieve the timestamp of the last input. This method has been detected on custom Delphi packers (FireEye, 2018). The last input is measured across multiple Human Interface Devices (HID) such as mice and keyboard. If the timestamp is old, the malware can say that no user is present on the system. The common Windows API call to get this information is `GetLastInputInfo` which returns the timestamp of the last input action in system ticks. Each tick is roughly 1 ms and the timestamp is zero when the system boots. Our tool logs the call and sets the returned

⁴ <https://www.microsoft.com/it-it/download/details.aspx?id=3702>

value of `GetLastInputInfo` to the following value: `LastInput = CurrentTicks() + rand()%1000`

A9. Registry

The Windows Registry is a container of valuable system information. It contains interesting values such as Services available, Programs Installed and System information that can also be retrieved using Windows API calls. The Registry is a hierarchical database where each key stores some values. Each value has some typed data associated that can be queried using Windows Registry Editor or programmatically using Windows Libraries. The type of detection is similar to the one described for File System Artifacts (Instinct, 2019). Indeed, we use this blacklist to log and then block suspicious Registry Keys. In particular, we focus on the most reliable techniques to retrieve Registry information.

NtOpenKey/Ex The only two syscalls available to open a registry key. The existence of a VM/Debugger related key is a good indicator of a malware analysis system. There are some VirtualBox specific Keys that are never present on physical machines. Our tool allows the generation of a valid HANDLE if the full key path does not contain blacklisted keywords, otherwise STATUS_OBJECT_NAME_NOT_FOUND is returned. An example of blacklisted key is: `HKLM\HARDWARE\ACPI\DSDT\VBOX_`.

NtQueryValueKey It takes the HANDLE of an open key and a value name and returns the data stored in the registry. If the value name contains some blacklisted keywords, the result is a non-existing value error. If the data stored contains a blacklisted keyword, the whole data is replaced with a random string. This implementation takes into account String (REG_SZ, REG_EXPAND_SZ) and MultiStrings (REG_MULTI_SZ) values.

NtEnumerateKey It enumerates all the subkeys of a key. Again, if some subkey should be blacklisted, its name is replaced with random values to hide its presence.

NtEnumerateValueKey It is used to enumerate all the values and data contained in a key. In this case, the same countermeasure of NtEnumerateKey is applied, taking into account REG_SZ, REG_EXPAND_SZ, REG_MULTI_SZ values.

A10. System environment

Physical Memory Information. Malware can detect analysis VMs by checking memory size information, because they are usually configured to take few GB of memory. An evasive behavior consists in checking if the physical memory size is significantly less than 4 GB (Project, 2017b). There are two ways to obtain this information: First, using `NtQuerySystemInformation(PHYSICAL_MEMORY_INFO)`, an undocumented syscall. Second, using `GlobalMemoryStatusEx`, a Windows API call that returns a documented and detailed report on the memory status, including free memory and total memory size.

Disk Space Information. The size of the disk installed in the system is a good indicator of VMs, because most VMs do not take more than 80 GB (Project, 2017b). There are three documented ways to obtain information about the disk size. `GetDiskFreeSpaceEx/A/W` returns the total space on disk and free bytes; `GetDiskFreeSpace/A/W` returns the same information and additional specifics on Sectors and Clusters sizes; `DeviceIoControl(IOCTL_DISK_GET_LENGTH_INFO)` is used to send control codes to device drivers. We use `IOCTL_DISK_GET_LENGTH_INFO` to ask for Disk Information.

The detection is done using the Windows API Hooker module, and the countermeasures applies a Disk Space Multiplier on the total size of the disk. We set this multiplier to 5, e.g., if the analysis VM has 30 GB, the API call will return 150 GB.

Network Information. It is widely known by malware developers that MAC Addresses have specific starting bytes depending on the Hypervisor Brand and the same Virtual NIC names are used across multiple VMs configurations. For example VirtualBox uses these starting bytes 08:00:27 for every NIC installed on the Guest OS (Cyberbit, 2007). An evasive malware simply can retrieve the MAC address of any NIC installed and perform a comparison on the first 3 bytes. NIC information can be retrieved by either calling `GetAdaptersInfo` or `GetAdaptersAddresses`. The two aforementioned Windows API calls are Hooked, Logged and the returned MAC Addresses and NIC Names are set to values defined by the analyst. For our analysis we used a MAC address corresponding to Intel Corporate Adapters and a NIC name that resembles the name of a real network card.

Computer Name detection. The computer name information is used to fingerprint sandboxes. To evade specific sandboxes, malware developers build a list of sandboxes computer names and then they use the functions `GetComputerName/Ex/A/W` to retrieve this information. Our tool intercepts attempts to retrieve this information.

IsDebuggerPresent . A Windows API that returns True if a debugger is attached to the process (Lindorfer et al., 2011; Sikorski and Honig, 2012).

CheckRemoteDebuggerPresent . This checks if the process is under a debugger in separate and parallel process (Lindorfer et al., 2011; Sikorski and Honig, 2012).

NtQuerySystemInformation(0x23) . Syscall that returns a structure containing a flag set to true when a Kernel Debugger is in place (Sikorski and Honig, 2012).

NtQueryObject . This technique makes use of undocumented system calls. The program begins creating a debug object using `NtCreateDebugObject`, and then the handle to that object is queried to see how many debug objects are present. In principle, there should be a debug object for each debugged process, however if a debug object is created and a debugger is present, two debug objects should be present. This is a very powerful feature and can be exploited also for anti-anti-debug checks. Our detection is to log the action if a `NtQueryObject` syscall is executed after a `NtCreateDebugObject` (Project, 2017a).

FindWindow . This technique is commonly used to detect Debuggers and Analyzers, and returns an Handle to the window corresponding to the window name specified as a parameter. Using this technique it is possible to know if specific tools are running on the system (Sikorski and Honig, 2012).

Number of Processors. Nowadays most machines have more than 1 CPU, however it is not uncommon to find VM configurations with a single CPU. Thus, the number of processors has been exploited as a feature by evasive malware to detect VM presence (McAfee, 2019). This information can be retrieved using Windows API calls, but the way most API calls get this value is reading a field inside the Process Environment Block (PEB) inside the Windows process address space. The number of CPU is at the `PEB_BASE_ADDRESS + 0x64`. Our tool using the Memory Accesses Controller monitors any access to this memory location logging the access and setting the value to 4. This way, any malware will detect a 4 CPU machine.

A11. WMI

Windows Management Instrumentation is a proprietary technology by Microsoft for Enterprise Management of Windows Machine. It is possible to write scripts for machines maintenance and retrieve their configurations. This last aspect can be exploited for evasion purposes. Indeed, WMI framework allows C++ developers and system administrators to write WMI Queries in WQL which is a SQL-like language for WMI information. These queries will re-

turn any kind of System related information. For example, performing `SELECT * FROM Win32_BIOS` it is possible to retrieve information about the BIOS Manufacturer which is specific for each Hypervisor. Other information can be extracted for fingerprinting purposes such as: process list, system resources, AVs installed and more (FireEye, 2016). The main method to execute a query is `ExecQuery` which requires the query string in the WQL language. Our bypass methodology is based on denying and logging every WMI query attempt. Indeed, it is unfeasible to understand from WQL queries which specific information the malware is searching.

A12. Process environment

`NtQueryInformationProcess` is a system call to retrieve information about the running process. Three Process Information Classes (under the form of integers) can be exploited to detect a debugger. `NtQueryInformationProcess(0x07)` is the underlying system call corresponding to the `CheckRemoteDebuggerPresent` API call (Sikorski and Honig, 2012). `NtQueryInformationProcess(0x1e)` is undocumented but leaks the presence of a debugger with the same behaviour of (1) (Kulchytsky and Kukoba, 2019). `NtQueryInformationProcess(0x1f)` is undocumented but behaves as (1) and (2) (Kulchytsky and Kukoba, 2019).

`NtSetInformationThread(0x11)`. This is a system call used to prevent the thread from sending debugging events to the attached debugger (Ferrie, 2011).

`NtCreateThread(HIDE_FROM_DEBUGGER)`. This is a system call used to tamper with debugging creating a thread that cannot be followed by the debugger (Kulchytsky and Kukoba, 2019).

`NtGetContextThread(CONTEXT_DEBUG_REGISTERS)`. This is used to get the value of Hardware Debug Registers. When Hardware breakpoints are used, the returned values are always different from zero because breakpoint addresses are set (Kulchytsky and Kukoba, 2019).

`NtSetContextThread(CONTEXT_DEBUG_REGISTERS)`. This is used to set the value of Hardware Debug Registers. In order to tamper with debugging it is possible to set to 0 every Debug Register in order to miss breakpoints (Kulchytsky and Kukoba, 2019).

DebugActiveProcess on Parent Process. A process cannot be debugged by two debuggers, so if a debugger is in place a child process with debugging privilege should fail its attempt to debug the parent process. Using this principle, the running process creates a child process and then the child process calls `DebugActiveProcess` on the parent. Therefore, a failed attempt to debug the parent is a clear indicator of the debugger presence (Ferrie, 2011). To handle this case our tool supports forking. The monitoring activity is extended to the child process and when `DebugActiveProcess` is called on the parent process, the action is logged through the Windows API Hooker.

A13. File system

An evasive behavior is to check inside known directories if a resource is present or not. For example, in Windows VirtualBox components are usually located under `C:\Program Files\Oracle\VirtualBox Guest Additions\`. A successful attempt to open this directory is an indicator of VM presence (Instinct, 2019). This same concept is applied to devices and files. There are multiple ways to access File System resources: First, using `NtOpenFile`, a syscall to open files, directories or devices. On success it returns an handle. When a blacklisted resource is opened, our tool makes the open fail with the return value `STATUS_OBJECT_NAME_NOT_FOUND`. Second, `NtCreateFile` is a syscall designed to supersede

`NtOpenFile` and it has the same functionalities, but it can also create new files and directories. The same behavior as `NtOpenFile` is applied when a blacklisted resource is created or opened. Third, `NtQueryAttributesFile` returns information about creation time, last access and write time, file size and file name length. Depending on the required attributes, file information can be more detailed. This syscall fails when the file is non existent, so a stealthier approach is to call it to check file presence. Also in this case our tool returns `STATUS_OBJECT_NAME_NOT_FOUND` when a blacklisted resource is detected. Finally, `NtQueryDirectoryFileEx` is a system call used to list the files contained in a directory. Before returning the whole list of files, our tool checks if some of them contain blacklisted keywords. In the latter case, the resource name is replaced with a random string.

A14. List processes

An easy way to detect debuggers, VMs and generic analysis tools is the process enumeration method (Instinct, 2019). The list of the processes can be retrieved using two methods. First, by calling `NtQuerySystemInformation(SYSTEM_PROCESS_INFORMATION)`, where `SYSTEM_PROCESS_INFORMATION = 5`, which is usually called by some high-level Windows API functions such as `CreateToolhelp32Snapshot`. Second, using the `EnumProcesses` Windows API, followed by `GetModuleBaseName`. Our tool is able to log attempts to retrieve the process list substituting at runtime the process names in blacklist with random strings, thus, fooling malware.

A15. List services

Services Enumeration. When a Guest OS is installed, hypervisor-related artifacts (e.g., VirtualBox Guest Additions) could be needed to work properly. These components could run as services inside the OS. Malware enumerates services to search for possible VM components (Instinct, 2019). This behavior works as follows. The Service Manager is opened via `OpenSCManager` Windows API call with the service enumeration permission. The actual Service Enumeration is performed using `EnumServicesStatus`, which returns a list of Service Names in the form of Keys for the Services Manager Database and a human readable format. If a malware sample knows exactly what it is searching for it can call `GetServiceKeyName`, `GetServiceDisplayName` or `OpenService` after `OpenSCManager`. The first requires the human readable name of the service, the second one requires the name assigned by the Service Manager Database, while the last one is used to open an existing service based on its name. We log the aforementioned calls and, if `EnumServicesStatus` reveals a name contained in our blacklist, the service name is substituted with random Display and Key names, while when `OpenService` is called on a blacklisted service, the action fails. If the detection of running VM services is done using `GetServiceKeyName` or `GetServiceDisplayName`, the value returned is an empty string, as specified by the documentation in case of failure.

A16. Drivers information

Drivers Enumeration. The Drivers enumeration evasion technique follows the same principle of Services Enumeration. In order to work correctly most Guest OS contains drivers related to Virtual Peripherals built by the Hypervisor. These components can be enumerated either by direct access to the Windows Driver

folder or using a Windows API call. The first method can be handled blocking access to VM files to the malware as we will discuss in Section A.13. Some examples of drivers checked by malware is described in DeepInstinct (2019). The second method exploits EnumDeviceDriver API call. An initial call is done to EnumDeviceDriver, which returns an array of drivers addresses. For each of them a call is done to GetDeviceDriverBaseName in order to get the Driver name. The Anti-Evasion System intercepts EnumDeviceDrivers as an action and replace any VM-related driver name returned by GetDeviceDriverBaseName with a random string.

References

- Afianian, A., Niksefat, S., Sadeghiyan, B., Baptiste, D., 2019. Malware dynamic analysis evasion techniques: a survey. *ACM Comput. Surv.* 52 (6).
- Barbosa, G.N., Branco, R.R., 2014. Prevalent characteristics in modern malware. Black Hat USA, Las Vegas, NV.
- Blackthorne, J., Bulazel, A., Fasano, A., Biernat, P., Yener, B., 2016. AVLeak: fingerprinting antivirus emulators through black-box testing. 10th USENIX Workshop on Offensive Technologies (WOOT 16).
- Branco, R.R., Barbosa, G.N., Drimel, P., 2012. Scientific but not academical overview of malware anti-debugging, anti-disassembly and anti-vm technologies. Black Hat USA, Las Vegas, NV.
- Bulazel, A., Yener, B., 2017. A survey on automated dynamic malware analysis evasion and counter-evasion: PC, mobile, and web. In: Proceedings of the 1st Reversing and Offensive-Oriented Trends Symposium. Association for Computing Machinery, New York, NY, USA.
- Checkpoint, 2020. Evasions: CPU. <https://evasions.checkpoint.com/techniques/cpu.html>.
- Chen, P., Huygens, C., Desmet, L., Joosen, W., 2016. Advanced or not? A comparative study of the use of anti-debugging and anti-vm techniques in generic and targeted malware. In: Hoepman, J.-H., Katzenbeisser, S. (Eds.), *ICT Systems Security and Privacy Protection*. Springer International Publishing, Cham, pp. 323–336.
- Collberg, C.S., Thomborson, C., 2002. Watermarking, tamper-proofing, and obfuscation-tools for software protection. *IEEE Trans. Softw. Eng.* 28 (8), 735–746.
- Comparetti, P.M., Salvaneschi, G., Kirda, E., Kolbitsch, C., Kruegel, C., Zanero, S., 2010. Identifying dormant functionality in malware programs. In: 2010 IEEE Symposium on Security and Privacy, pp. 61–76.
- Continella, A., Carminati, M., Polino, M., Lanzi, A., Zanero, S., Maggi, F., 2017. Prometheus: analyzing Webinject-based information stealers. *J. Comput. Secur.* 25 (2), 117–137. doi:10.3233/JCS-15773.
- Continella, A., Guagnelli, A., Zingaro, G., Pasquale, G.D., Barengi, A., Zanero, S., Maggi, F., 2016. ShieldFS: a self-healing, ransomware-aware filesystem. In: Schwab, S., Robertson, W.K., Balzarotti, D. (Eds.), *Proceedings of the 32nd Annual Conference on Computer Security Applications, ACSAC 2016, Los Angeles, CA, USA, December 5-9, 2016*. ACM, pp. 336–347. <http://dl.acm.org/citation.cfm?id=2991110>
- Cyberbit, 2007. Anti-vm and anti-sandbox explained. <https://www.cyberbit.com/blog/endpoint-security/anti-vm-and-anti-sandbox-explained/>.
- DeepInstinct, 2019. Anti-virtualization malware. <https://www.deepinstinct.com/2019/10/29/malware-evasion-techniques-part-2-anti-vm-blog/>.
- D'Elia, D.C., Coppa, E., Nicchi, S., Palmaro, F., Cavallaro, L., 2019. SoK: using dynamic binary instrumentation for security (and how you may get caught red handed). In: *Proceedings of the 2019 ACM Asia Conference on Computer and Communications Security*. ACM, New York, NY, USA, pp. 15–27.
- D'Elia, D.C., Coppa, E., Palmaro, F., Cavallaro, L., 2020. On the dissection of evasive malware. *IEEE Trans. Inf. Forensics Secur.* 15, 2750–2765.
- Ferrie, P., 2011. The ultimate anti-debug reference. https://anti-reversing.com/Downloads/Anti-Reversing/The_Ultimate_Anti-Reversing_Reference.pdf.
- FireEye, 2014. Turing test in reverse: new sandbox-evasion techniques seek human interaction. <https://www.fireeye.com/blog/threat-research/2014/06/turing-test-in-reverse-new-sandbox-evasion-techniques-seek-human-interaction.html>.
- FireEye, 2016. Increased use of WMI for environment detection and evasion. https://www.fireeye.com/blog/threat-research/2016/10/increased_use_ofwmi.html.
- FireEye, 2018. Increased use of a delphi packer to evade malware classification. <https://www.fireeye.com/blog/threat-research/2018/09/increased-use-of-delphi-packer-to-evade-malware-classification.html>.
- Fratantonio, Y., Bianchi, A., Robertson, W., Kirda, E., Kruegel, C., Vigna, G., 2016. TriggerScope: towards detecting logic bombs in android applications. In: 2016 IEEE Symposium on Security and Privacy (SP). IEEE, pp. 377–396.
- Gagnon, M.N., Taylor, S., Ghosh, A.K., 2007. Software protection through anti-debugging. *IEEE Secur. Privacy* 5 (3), 82–84.
- Holz, R., Perino, D., Varvello, M., Amann, J., Continella, A., Evans, N., Leontiadis, I., Natoli, C., Scheitle, Q., 2020. A retrospective analysis of user exposure to (illicit) cryptocurrency mining on the web. In: *Proceedings of the Network Traffic Measurement and Analysis Conference (TMA)*.
- Instinct, D., 2017. Common anti-debugging techniques in the malware landscape. <https://www.deepinstinct.com/2017/12/27/common-anti-debugging-techniques-in-the-malware-landscape/>.
- Instinct, D., 2019. Anti-virtualization malware. <https://www.deepinstinct.com/2019/10/29/malware-evasion-techniques-part-2-anti-vm-blog/>.
- Kharraz, A., Ma, Z., Murley, P., Lever, C., Mason, J., Miller, A., Borisov, N., Antonakakis, M., Bailey, M., 2019. Outguard: detecting in-browser covert cryptocurrency mining in the wild. In: *Proceedings of the The World Wide Web Conference (WWW)*.
- Kharraz, A., Robertson, W., Balzarotti, D., Bilge, L., Kirda, E., 2015. Cutting the gordian knot: a look under the hood of ransomware attacks. In: *Proceedings of the International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment (DIMVA)*. Springer.
- Kirat, D., Vigna, G., Kruegel, C., 2014. Barecloud: bare-metal analysis-based evasive malware detection. In: *Proceedings of the 23rd USENIX Conference on Security Symposium*. USENIX Association, Berkeley, CA, USA, pp. 287–301.
- Kolbitsch, C., Kirda, E., Kruegel, C., 2011. The power of procrastination: detection and mitigation of execution-stalling malicious code. In: *Proceedings of the 18th ACM Conference on Computer and Communications Security*. ACM, New York, NY, USA, pp. 285–296.
- Konoth, R.K., Vineti, E., Moonsamy, V., Lindorfer, M., Kruegel, C., Bos, H., Vigna, G., 2018. Minesweeper: an in-depth look into drive-by cryptocurrency mining and its defense. In: *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security (CCS)*.
- Küchler, A., Mantovani, A., Han, Y., Bilge, L., Balzarotti, D., 2021. Does every second count? Time-based evolution of malware behavior in sandboxes. In: *Proceedings of the ISOC Network and Distributed System Security Symposium (NDSS)*.
- Kulchitsky, O., Kukoba, A., 2019. Anti debugging protection techniques with examples. <https://www.apriorit.com/dev-blog/367-anti-reverse-engineering-protection-techniques-to-use-before-releasing-software>.
- Labs, M., 2017. McAfee labs threats report. <https://www.mcafee.com/enterprise/en-us/assets/reports/rp-quarterly-threats-jun-2017.pdf>.
- Lindorfer, M., Kolbitsch, C., Milani Comparetti, P., 2011. Detecting environment-sensitive malware. In: Sommer, R., Balzarotti, D., Maier, G. (Eds.), *Recent Advances in Intrusion Detection*. Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 338–357.
- McAfee, 2017. McAfee Labs Threats Report. <https://www.mcafee.com/enterprise/en-us/assets/reports/rp-quarterly-threats-jun-2017.pdf>.
- McAfee, 2019. Evolution of malware sandbox evasion tactics a retrospective study. <https://securingtomorrow.mcafee.com/other-blogs/mcafee-labs/evolution-of-malware-sandbox-evasion-tactics-a-retrospective-study/>.
- McAfee.com, 2018. McAfee labs threats report. <https://www.mcafee.com/enterprise/en-us/assets/reports/rp-quarterly-threats-dec-2018.pdf>.
- Microsoft, 2020. Determining if a debugger is attached. <https://docs.microsoft.com/en-us/windows-hardware/drivers/devtest/determining-if-a-debugger-is-attached>.
- Miramirkhani, N., Appini, M.P., Nikiforakis, N., Polychronakis, M., 2017. Spotless sandboxes: evading malware analysis systems using wear-and-tear artifacts. In: *Proceedings of the IEEE Symposium on Security and Privacy (S&P)*.
- Oyama, Y., 2018. Investigation of the diverse sleep behavior of malware. *J. Inf. Process.* 26, 461–476.
- Oyama, Y., 2018. Trends of anti-analysis operations of malwares observed in API call logs. *J. Comput. Virol. Hacking Tech.* 14 (1), 69–85.
- Paleari, R., Martignoni, L., Roglia, G.F., Bruschi, D., 2009. A fistful of red-pills: how to automatically generate procedures to detect cpu emulators. In: *Proceedings of the USENIX Workshop on Offensive Technologies (WOOT)*, vol. 41, p. 86.
- Paranoid Fish (Pafish) <https://github.com/aOrtega/pafish>. 2020 PEID2020. <https://www.aldeid.com/wiki/PEID>.
- Polino, M., Continella, A., Mariani, S., D'Alessio, S., Fontana, L., Gritti, F., Zanero, S., 2017. Measuring and defeating anti-instrumentation-equipped malware. In: Polychronakis, M., Meier, M. (Eds.), *Detection of Intrusions and Malware, and Vulnerability Assessment*. Springer International Publishing, Cham, pp. 73–96.
- Project, U., 2017a. Anti-debugging. <http://unprotect.tdgt.org/index.php/Anti-debugging>.
- Project, U., 2017b. Sandbox evasion. http://unprotect.tdgt.org/index.php/Unprotect_Project.
- Quarta, D., Salvioni, F., Continella, A., Zanero, S., 2018. Toward systematically exploring antivirus engines. In: *Proceedings of the Conference on Detection of Intrusions and Malware and Vulnerability Assessment (DIMVA)*.
- Raffetseder, T., Kruegel, C., Kirda, E., 2007. Detecting system emulators. In: *International Conference on Information Security*. Springer, pp. 1–18.
- Rossov, C., Dietrich, C.J., Grier, C., Kreibich, C., Paxson, V., Pohlmann, N., Bos, H., Van Steen, M., 2012. Prudent practices for designing malware experiments: status quo and outlook. In: 2012 IEEE Symposium on Security and Privacy. IEEE, pp. 65–79.
- Rutkowska, J., 2004. Redpill: detect VMM using (almost) one cpu instruction. <http://invisiblethings.org/papers/redpill.html>.
- Sebastián, M., Rivera, R., Kotzias, P., Caballero, J., 2016. AVCLASS: a tool for massive malware labeling. In: *International Symposium on Research in Attacks, Intrusions, and Defenses*. Springer, pp. 230–253.
- Sikorski, M., Honig, A., 2012. *Practical Malware Analysis: The Hands-On Guide to Dissecting Malicious Software*, 1st ed. No Starch Press, San Francisco, CA, USA.
- Timzen, T., 2020. Kernel Forensics and Rootkits. https://www.tophertimzen.com/resources/cs407/slides/week06_01-Rootkits.html.
- Ugarte-Pedrero, X., Balzarotti, D., Santos, I., Bringas, P.G., 2015. SoK: deep packer inspection: a longitudinal study of the complexity of run-time packers. In: *Proc. of IEEE symposium on Security and Privacy (SP)*. IEEE.
- Virtual machine detection techniques. 2017. <https://shasaurabh.blogspot.com/2017/07/virtual-machine-detection-techniques.html>.
- Xu Chen, Andersen, J., Mao, Z.M., Bailey, M., Nazario, J., 2008. Towards an understanding of anti-virtualization and anti-debugging behavior in modern malware.

In: 2008 IEEE International Conference on Dependable Systems and Networks With FTCS and DCC (DSN), pp. 177–186.

Zhu, S., Shi, J., Yang, L., Qin, B., Zhang, Z., Song, L., Wang, G., 2020. Measuring and modeling the label dynamics of online anti-malware engines. USENIX Security Symposium).

Nicola Galloro received a M.Sc. degree (2017) in Computer Science and Engineering from Politecnico di Milano in Italy, where he worked as a researcher on computer security topics as part of the cybersecurity research group. His research was focused on evasive malware targeting Windows platforms. Currently, he is working as a software engineer, and before working as a researcher, he was a security consultant.

Mario has got his PhD (2017) in Computer Security from Politecnico di Milano. Currently, Mario is a PostDoc Security Research at Politecnico di Milano. His PhD was mainly based on Malware Analysis and automation of behavior extraction. Besides malware analysis, Mario is interested in several aspects of Computer Security; his works range from the study of Cyber-Physical systems to Banking Fraud Detection going through Android Security. Even more, he is interested in Binary analysis and Exploitation techniques. Mario spends most of his free time playing Capture the Flag Competition with "Tower of Hanoi" and "mhackeroni".

Michele Carminati received a Ph.D. degree in Information Technology from Politecnico di Milano in Italy, where he is currently a Postdoctoral Researcher working as part of the cybersecurity group. His research focuses on the application of machine learning methods in various security-related fields, ranging from cyber-physical and

automotive systems to binary analysis, going through anomaly and intrusion detection. In particular, he is actively involved in the financial fraud detection domain, where he works on the analysis of advanced financial threats such as banking Trojans and on the design of novel frameworks to identify and timely detect fraudulent transactions.

Andrea Continella is an Assistant Professor at the Faculty of Electrical Engineering, Mathematics and Computer Science of the University of Twente. Previously, he was a Postdoctoral Researcher in the Computer Science Department at UC Santa Barbara, and he obtained a Ph.D. cum laude in Computer Science and Engineering at Politecnico di Milano. His research activity focuses on different aspects of systems security, such as malware analysis, mobile and IoT security & privacy, automated vulnerability discovery, and large-scale measurements of security issues. Andrea has published several research papers and he serves in the program committees of well-known systems security conferences and workshops. Andrea also loves Capture The Flag (CTF) competitions, which he currently plays with Shellphish, and he co-organized several editions of the PoliCTF and iCTF.

Stefano Zanero received a PhD in Computer Engineering from Politecnico di Milano, where he is currently an associate professor with the Dipartimento di Elettrotecnica, Informazione e Bioingegneria. His research focuses on malware analysis, cyberphysical security, and cybersecurity in general. He is a Senior Member of the IEEE and sits in the Board of Governors of the IEEE Computer Society; he is a lifetime senior member of the ACM and has been named a Fellow of ISSA (Information System Security Association).