

Adversarial EXEmples: Evading Windows Malware Detectors



Luca Demetrio

luca.demetrio@unige.it

X @zangobot

Part 1 - Malware and how to detect it

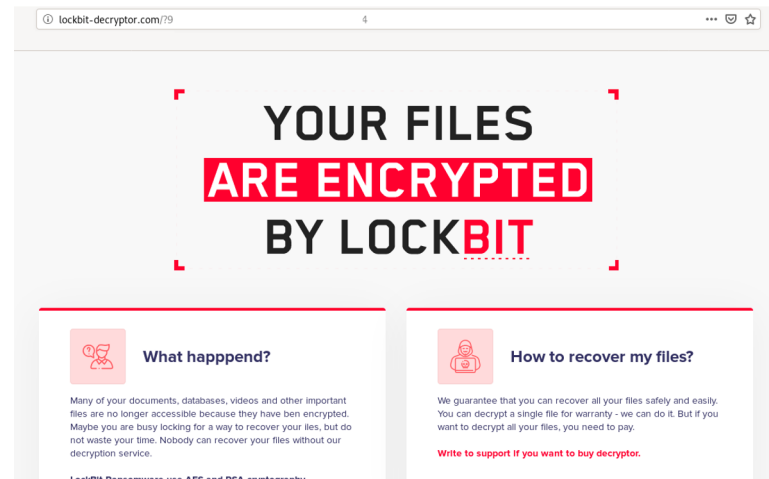
Malware: dangerous programs

Malicious by design

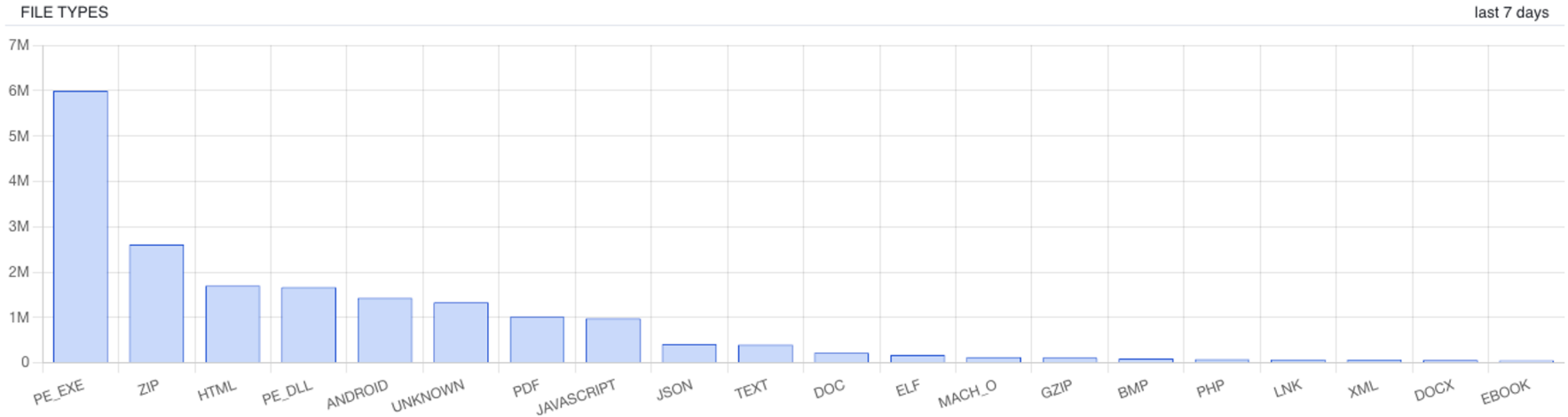
Intrusive programs whose intent is causing damage, steal information, ransom, take control of devices

“Business” of malware

Currently, there are plenty of famous active groups that act as companies that produce and sell malware as a service



Concerning numbers

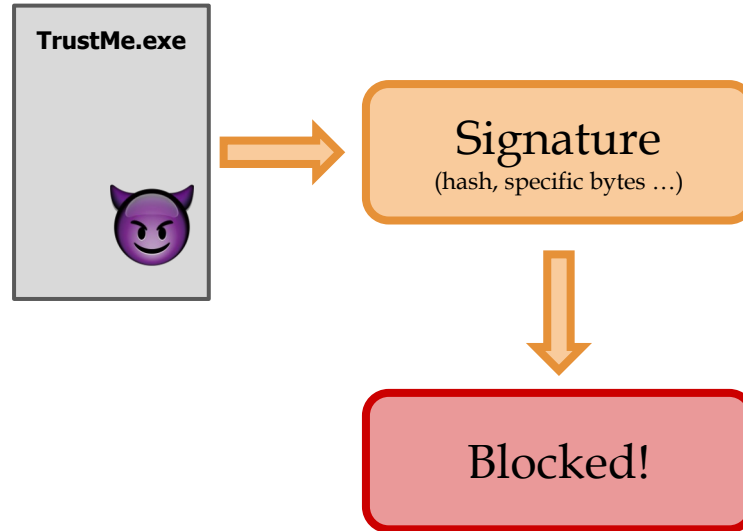


Every week, **6M** of Windows executable are scanned with cloud systems to detect malware, with more than **1M** detected samples in the same timespan

Security *without* Machine Learning

Blocklist approach

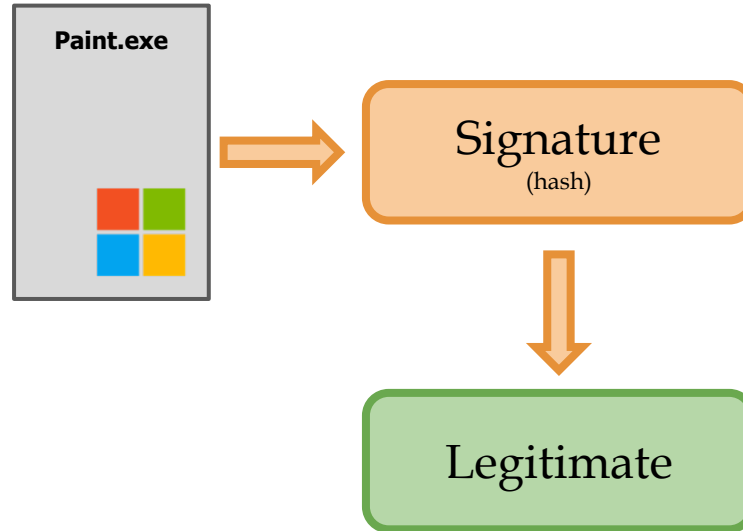
1. Extract a “**signature**” from each incoming request, like hashes, the presence of specific bytes or word (like “Viagra” for spam)
2. If these signatures are contained inside well-known **blocklists**, the input is recognized as **malicious**



Security *without* Machine Learning

Allowlist approach

1. Extract a “**signature**” from each incoming request, like hashes.
2. If these signatures are contained inside well-known **allowlists**, the input is recognized as **harmless**, otherwise it is blocked



YARA Rules

Known patterns as text

A rule is a sequence of bytes, and files match depending on the specified condition

Plenty of open-source rules

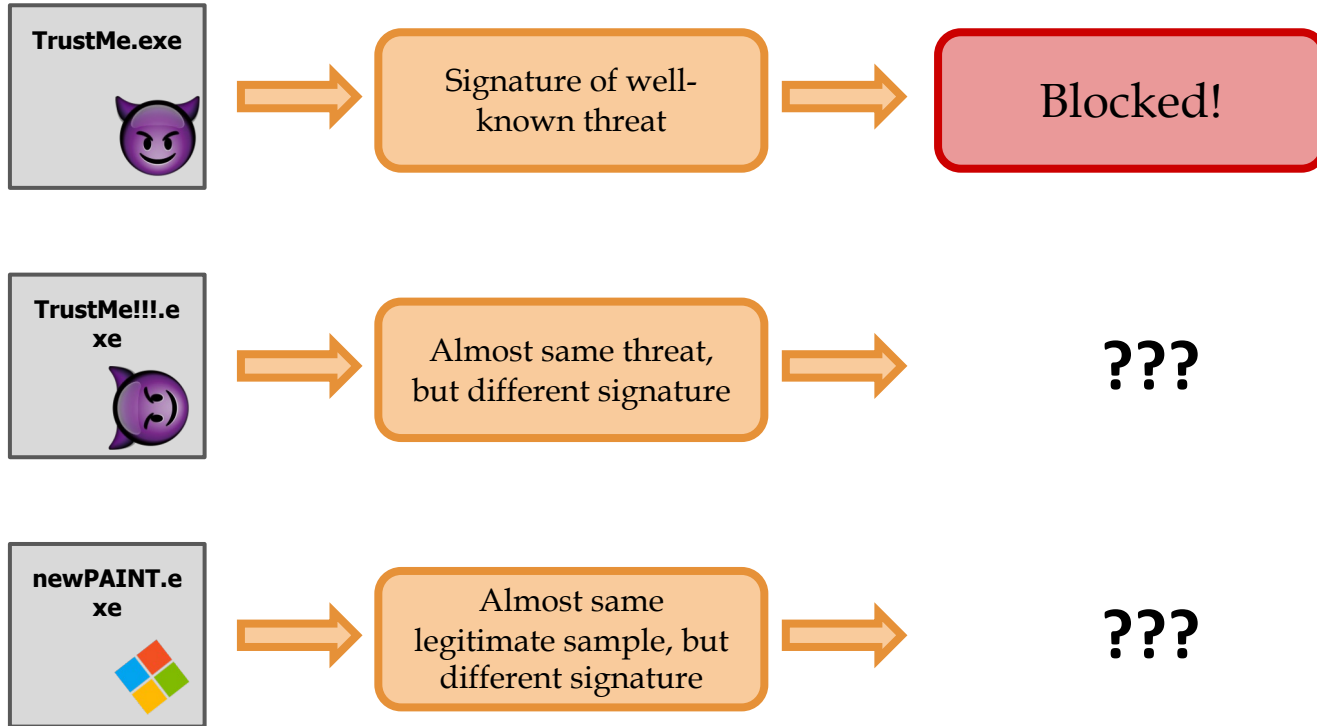
Try to google “YARA rules” :D

```
rule silent_banker : banker
{
    meta:
        description = "This is just an example"
        threat_level = 3
        in_the_wild = true

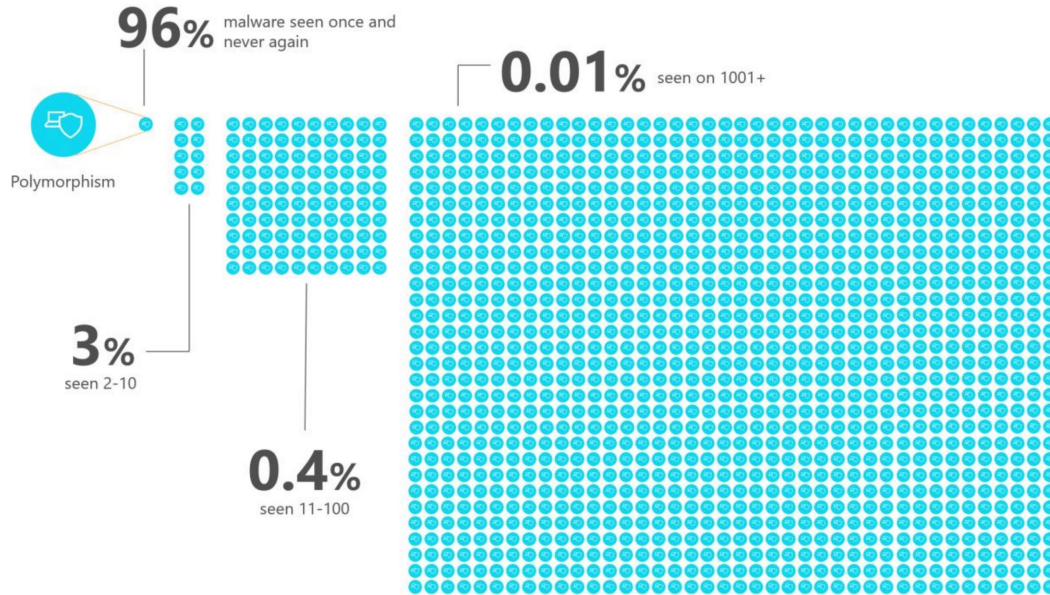
    strings:
        $a = {6A 40 68 00 30 00 00 6A 14 8D 91}
        $b = {8D 4D B0 2B C1 83 C0 27 99 6A 4E 59 F7 F9}
        $c = "UVODFRYSIHLNWPEJXQZAKCBGMT"

    condition:
        $a or $b or $c
}
```

Limitation: what happens with minimal changes?



Required amount of static signatures



Too many variants!

Block/Allowlist can work ONLY if threats are known, but the majority of them (96%) appear in the wild only **ONE** time

Suboptimal performances

Quality of rules matters

Open-source rules are not the best around, and it is possible that companies possess better ones for detection

(Thanks to Andrea Ponte, Ph.D. @ SMARTLAB for the initial analysis!)

```
Compiled rules:2563
Tested malwares:5005
Detected malwares: 765
Total time taken: 69.11032915115356
Average time taken: 0.013808257572658056
Standard deviation: 0.038165608714924255
```

Machine Learning for Malware Detection

Machine Learning applied as Security Scanner



Spreading into commercial products

Companies claim to use machine learning technologies inside their detectors to spot Windows malware by learning patterns from data

Filter out known threats, generalize to variants

Deep networks learn “signatures”, and they can spot variants of the same malware

Intercept X Deep Learning



CrowdStrike Introduces Enhanced Endpoint Machine Learning Capabilities and Advanced Endpoint Protection Modules

— Company continues to accelerate pace of replacement of legacy AV solutions in both enterprise and SMB markets —



Solutions ▾ Industries ▾ Products ▾ Services ▾ Resource Center ▾ About Us ▾ GDPR

Home > Enterprise Security > WMI > Machine Learning in Cybersecurity

TECHNOLOGY
Machine Learning in Cybersecurity

Static malware detection with ML

Programs stored as file

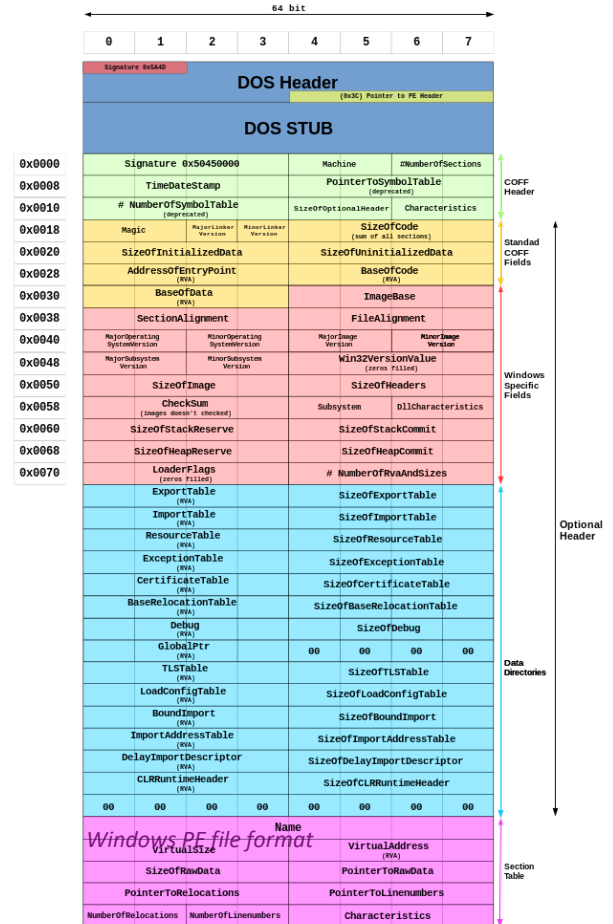
Each program is a regular file that can be analysed without executing it

Features to extract

Which API they import? How many resources they contain? Are there some initialized values? Do they require special permissions?

End-to-end learning

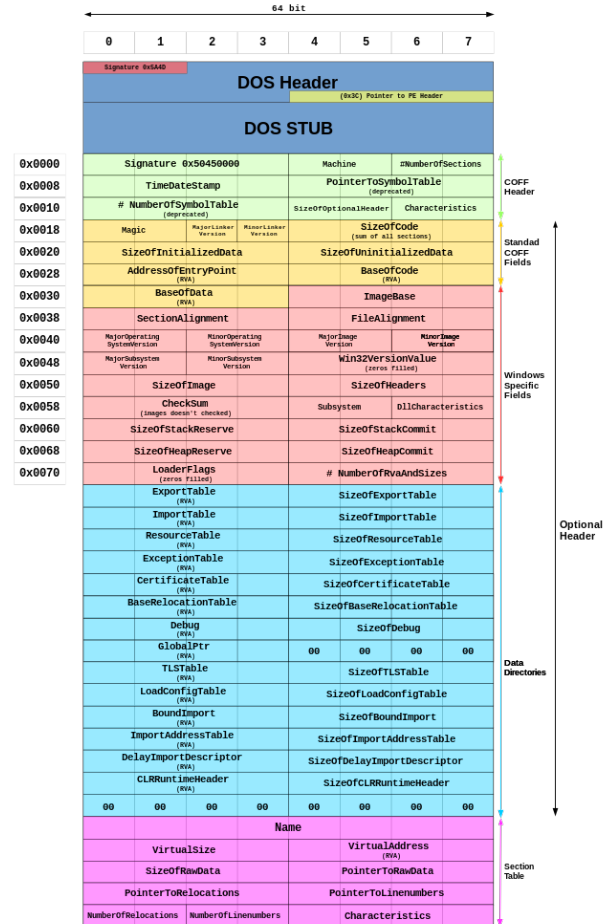
Or, use each byte of the program/resources/files as token, let the network learn by itself a suitable representation



Windows PE File Format

Format adapted for “modern” programs
(from Windows NT 3.1 on)

Before there were other formats, one is the
DOS (kept for retrocompatibility)

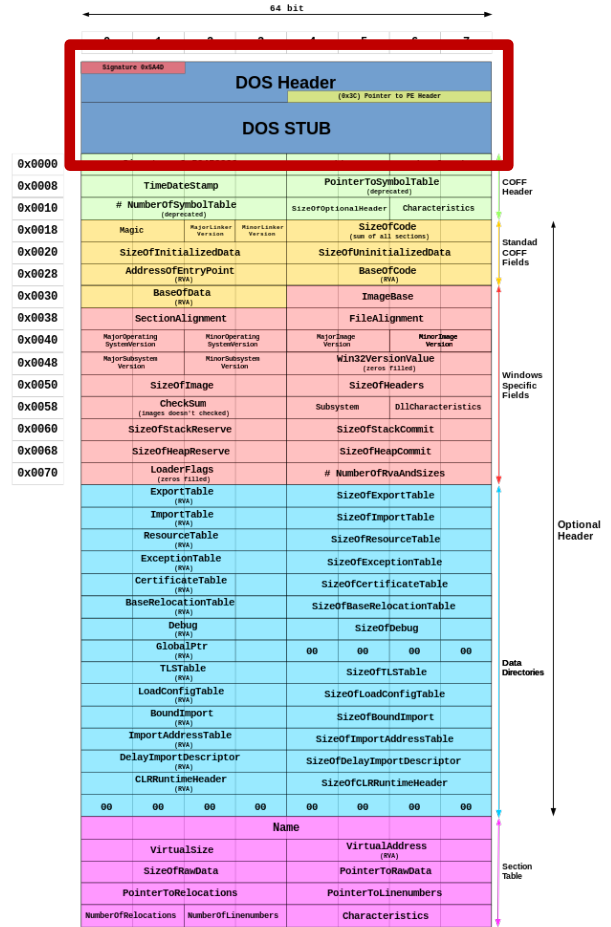


Windows PE File Format

DOS Header + Stub

Metadata for DOS program

Executing a modern program in DOS will trigger the “This program cannot be run in DOS mode” output

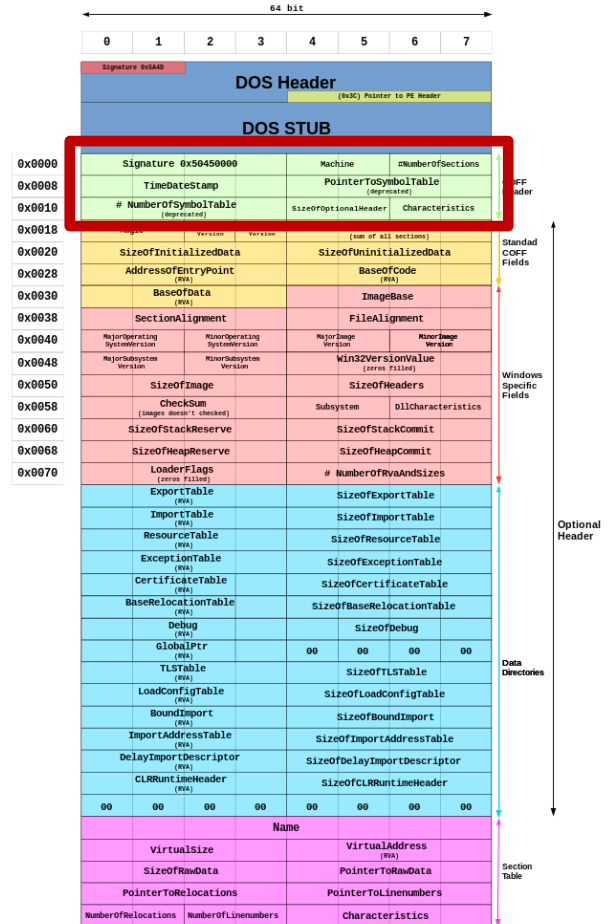


Windows PE File Format

PE Header

Real metadata of the program

Describes general information of the file

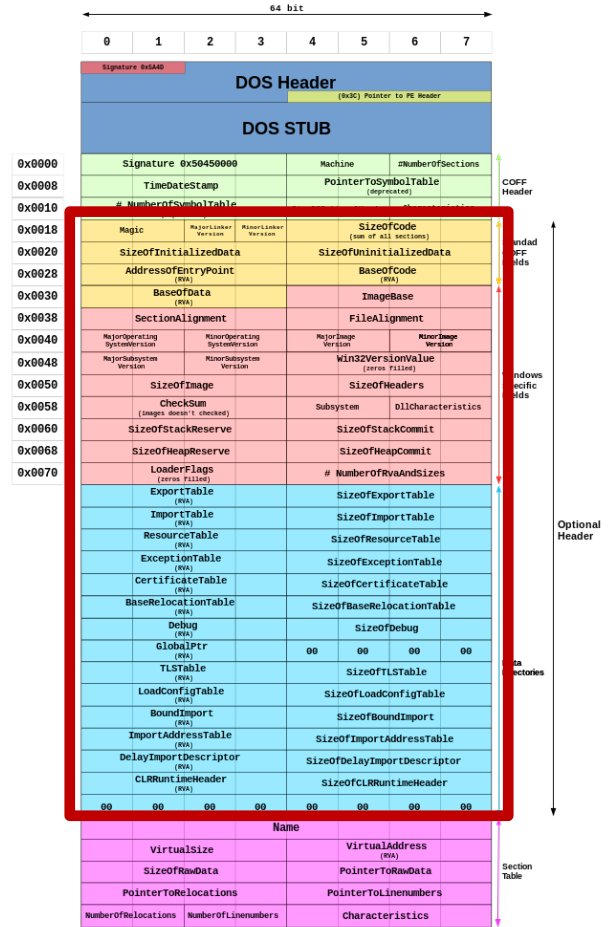


Windows PE File Format

Optional Header

Spoiler: not optional at all :)

Instructs the loader where to find each object inside the file



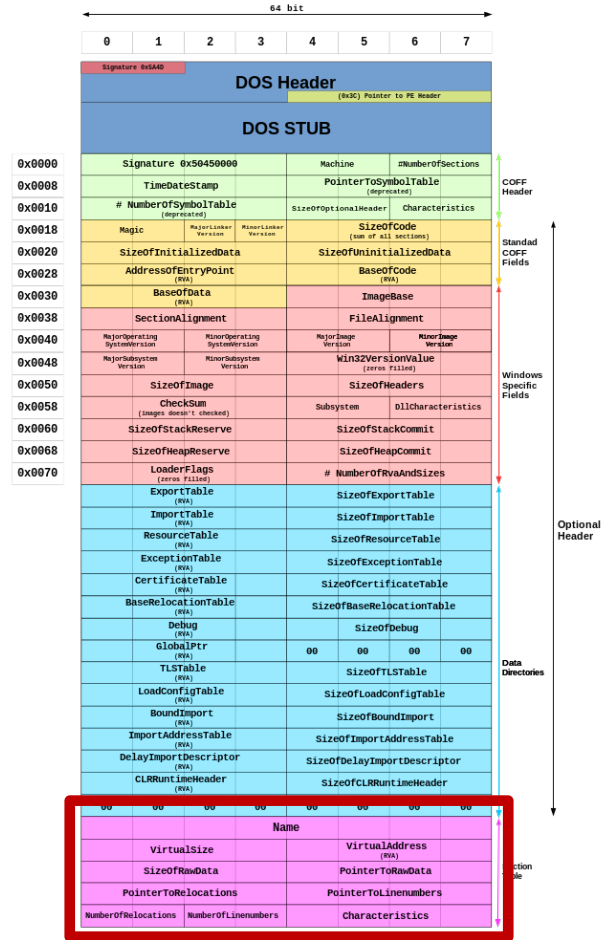
Windows PE File Format

Section Table and Sections

Describes where to find code, initialized data, resources, etc to the loader

These are “sections”, and each has a “section entry” with its characteristics

Examples: code is “.text”, read-only data is “.rodata”, resources are “.rsc”, and counting



How programs are loaded

1 Headers

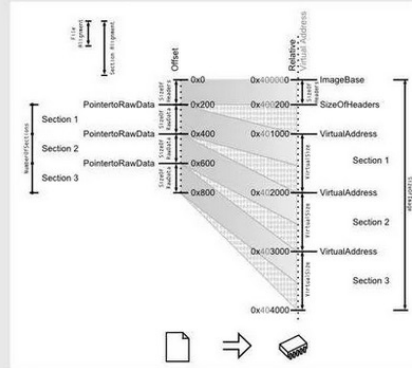
the *DOS Header* is parsed
the *PE Header* is parsed
(its offset is *DOS Header's e_lfanew*)
the *Optional Header* is parsed
(it follows the *PE Header*)

2 Sections table

Sections table is parsed
(it is located at: $\text{offset}(\text{OptionalHeader}) + \text{SizeOfOptionalHeader}$)
it contains *NumberOfSections* elements
it is checked for validity with alignments:
FileAlignments and *SectionAlignments*

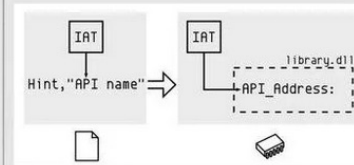
3 Mapping

the file is mapped in memory according to:
the *ImageBase*
the *SizeOfHeaders*
the *Sections* table



4 Imports

DataDirectories are parsed
they follow the *OptionalHeader*
their number is *NumOfRVAAndSizes*
imports are always #2
Imports are parsed
each descriptor specifies a *DllName*
this DLL is loaded in memory
IAT and *INT* are parsed simultaneously
for each API in *INT*
its address is written in the *IAT* entry



5 Execution

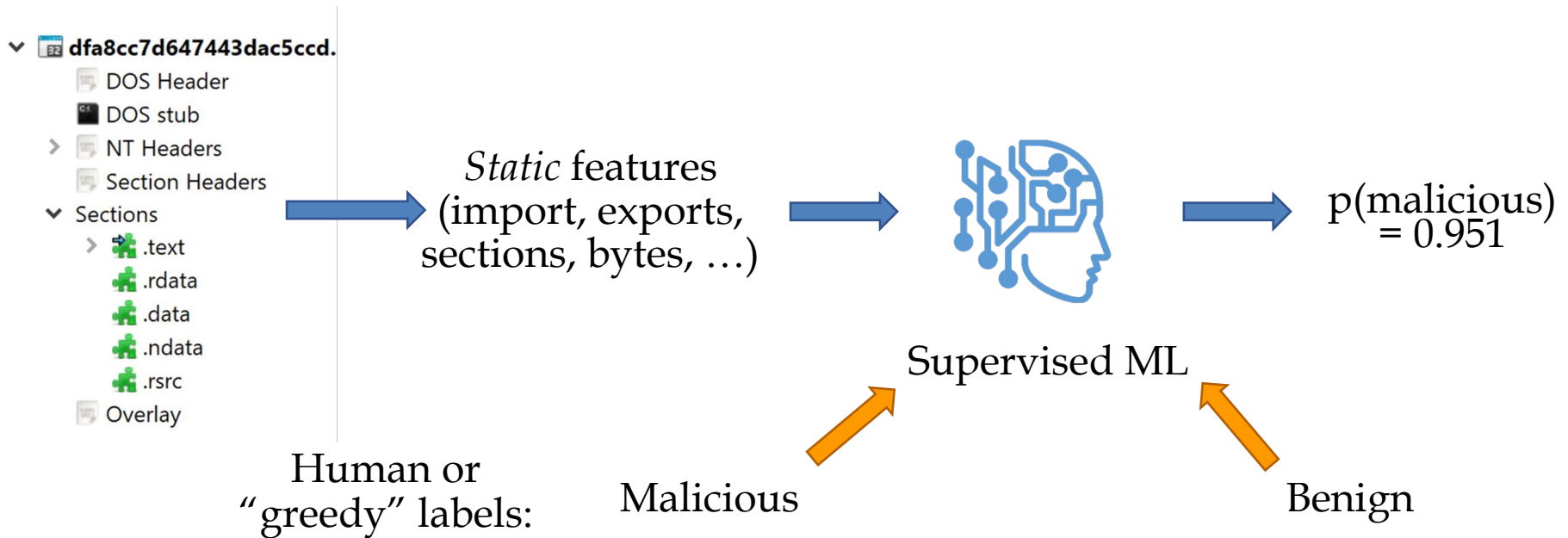
Code is called at the *EntryPoint*
the calls of the code go via the *IAT* to the APIs



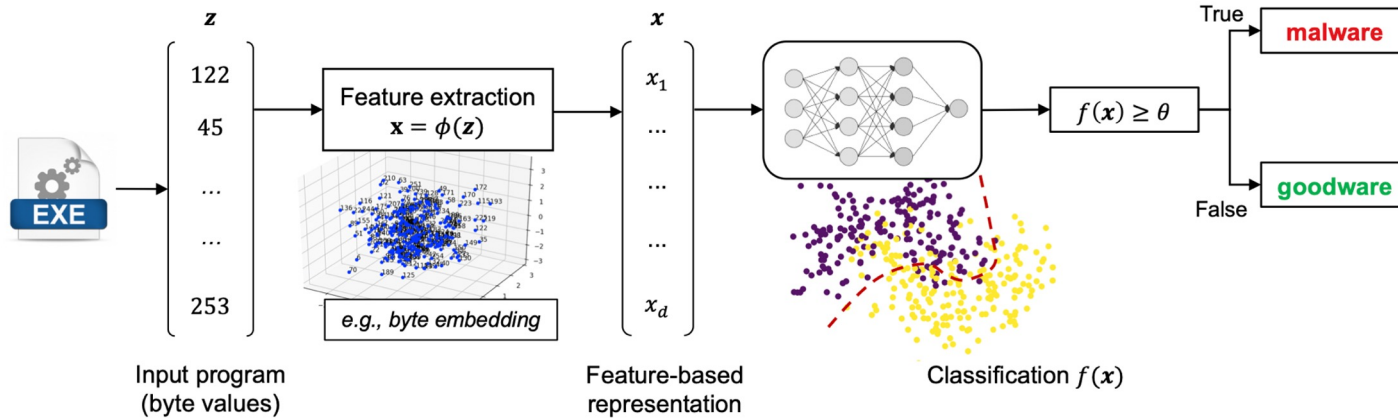
<https://code.google.com/archive/p/corkami/wikis/PE101.wiki>

(Thanks Ange Albertini)

Static malware detection with ML



Examples of Static Windows Malware Detectors



Raw bytes as features

MalConv => **1 MB** in input, embedding

Static features from data

GBDT => decision tree on **2.381** features
(API, sections, byte entropy, exports, strings, etc.)

MalConv: end-to-end classification

Total params to train: **1.042.953**

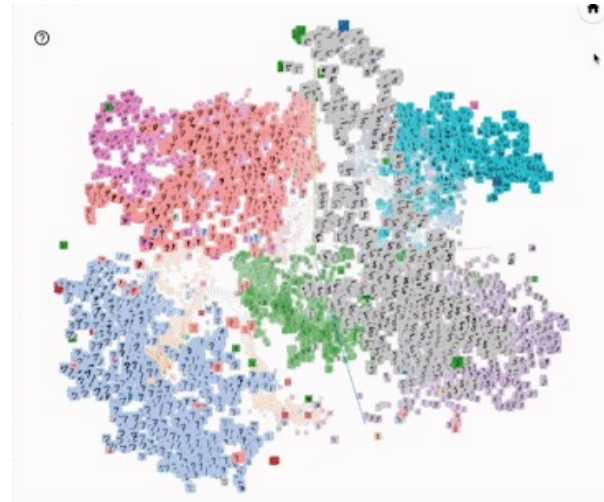
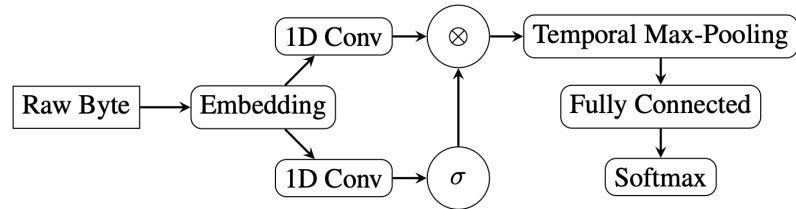
Embedding layer

Each byte is converted in a 8-dimensional vector learned at training time.

Convolutions similar to images

Analyse 1D patterns of bytes to retrieve local information, later aggregated by the fully-connected layer

(Thanks to Dmitrijs Trizna for the animation!)



GBDT EMBER: classic feature extraction

Total params to train: **2381**

Hand-crafted features

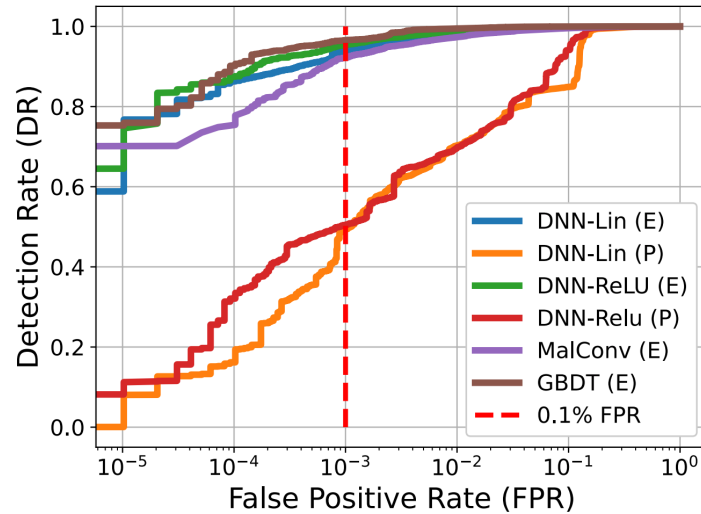
Sections, histogram of bytes, distribution of strings, API calls, etc

Harder to develop, faster to test

Instilling domain knowledge is hard, but the gradient-boosting decision tree model is fast to train (faster than a neural network)



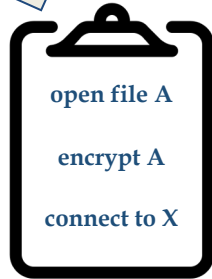
Performance comparison



Hand-crafted features seem better

Computed on the EMBER dataset, GBDT exhibits superior performance to other models

Dynamic malware detection with ML



Chain of events

Run program inside protected isolated environment, take note of every observable action of the program

Human-readable reports

The analysis outputs a textual report that specifies the timeline of all the triggered events

Circumventing obfuscation

Even if samples are packed or obfuscated, at some point the functionality will be manifested through interactions with the underlying OS

Dynamic malware detection with ML

dfa8cc7d647443dac5ccd.

- DOS Header
- DOS stub
- NT Headers
- Section Headers
- Sections
 - .text
 - .rdata
 - .data
 - .ndata
 - .rsrc
 - Overlay

→ **Sandbox** →

→ **Emulator** →

The screenshot displays the Process Hacker interface. The top bar shows 'Summary | 103 calls | 46 KB used | regsvr32.exe'. The main window is divided into several panes:

- Process Hacker:** Shows a tree view of processes. 'explorer.exe' is expanded, showing 'cmd.exe' and 'powershell.exe'. 'powershell.exe' is further expanded to show 'PDFXCview.exe'.
- Module:** Lists loaded modules for 'regsvr32.exe', including 'KERNELBASE.dll' and 'regsvr32.exe'.
- API:** A table of API calls made by 'regsvr32.exe'. The 'LdrLoadDll' and 'GetProcAddress' calls are highlighted with a red box. The 'GetProcAddress' calls show the function name 'ZwQuerySystemInformation'.
- Process Details:** Shows the command line for 'PDFXCview.exe' as '"C:\ProgramData\PDFXCview.exe"'. Other details include the current directory 'C:\ProgramData\' and the started time 'a minute and 32 seconds ago (8:18:57 A)'. The PEB address is '0x214000 (32-bit: 0x215000)'.
- Process List:** A table at the bottom shows the process name, PID, operation, and path for 'PDFXCview.exe' instances. The 'Process Create' operation is highlighted with a red box.

Dynamic malware detection with ML

Traces of execution

Structure alone can't reveal too much: exploit behavior of analysed programs

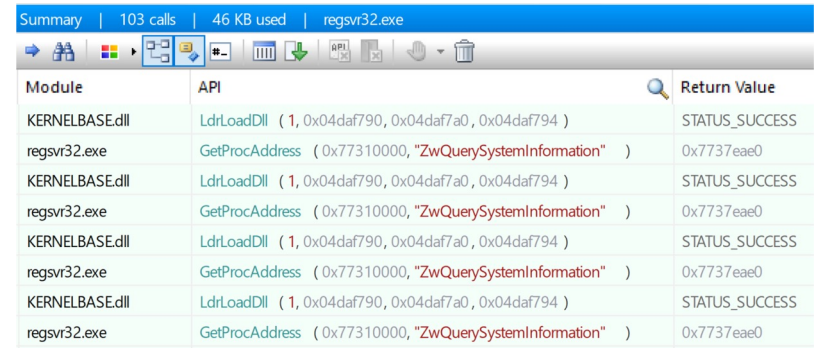
Running programs becomes necessary

Features to extract

Which API they call? Which IP addresses they contact? Which service they interact with?

End-to-end learning

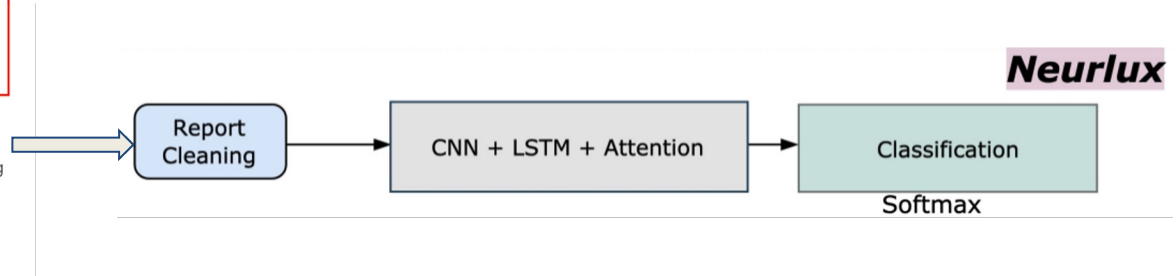
Or, consider the sequence of actions as a list of token that can be used in NLP systems



Module	API	Return Value
KERNELBASE.dll	LdrLoadDll (1, 0x04daf790, 0x04daf7a0, 0x04daf794)	STATUS_SUCCESS
regsvr32.exe	GetProcAddress (0x77310000, "ZwQuerySystemInformation")	0x7737eae0
KERNELBASE.dll	LdrLoadDll (1, 0x04daf790, 0x04daf7a0, 0x04daf794)	STATUS_SUCCESS
regsvr32.exe	GetProcAddress (0x77310000, "ZwQuerySystemInformation")	0x7737eae0
KERNELBASE.dll	LdrLoadDll (1, 0x04daf790, 0x04daf7a0, 0x04daf794)	STATUS_SUCCESS
regsvr32.exe	GetProcAddress (0x77310000, "ZwQuerySystemInformation")	0x7737eae0
KERNELBASE.dll	LdrLoadDll (1, 0x04daf790, 0x04daf7a0, 0x04daf794)	STATUS_SUCCESS
regsvr32.exe	GetProcAddress (0x77310000, "ZwQuerySystemInformation")	0x7737eae0

Example of Dynamic Windows Malware Detector

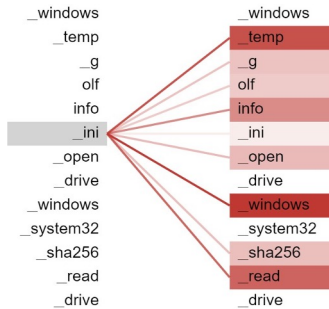
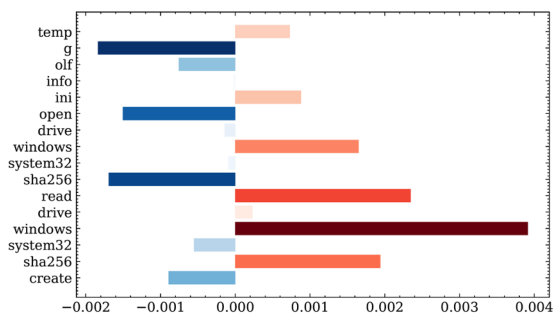
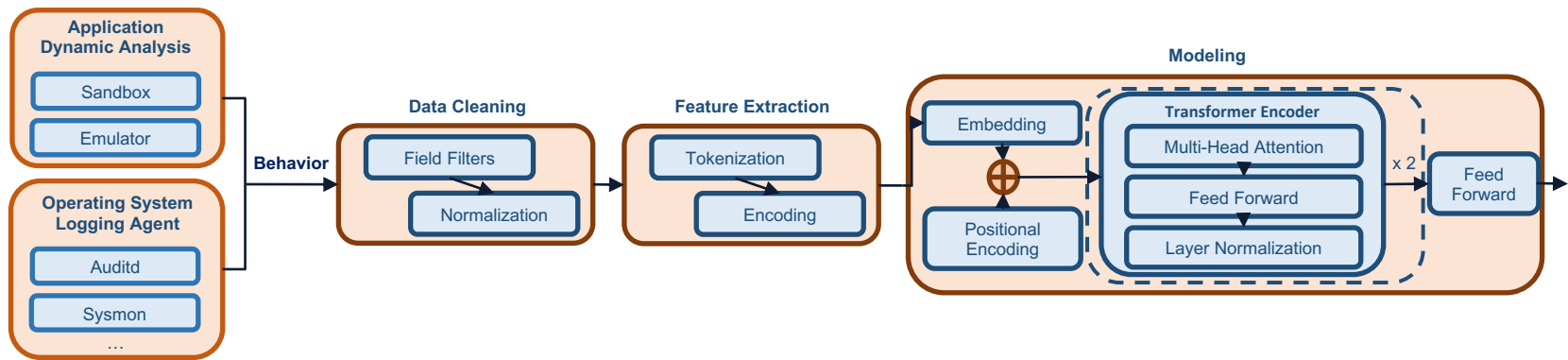
```
- HKLM\software\microsoft\windows nt\currentversion\winlogon\userinit  
  | C:\Windows\system32\userinit.exe,C:\Windows\apppatch\svchost.exe,  
+ HKLM\Software\Microsoft\Tracing\svchost_RASAPI32\EnableFileTracing  
+ HKLM\Software\Microsoft\Tracing\svchost_RASAPI32\EnableConsoleTracing  
+ HKLM\Software\Microsoft\Tracing\svchost_RASAPI32\FileTracingMask  
+ HKLM\Software\Microsoft\Tracing\svchost_RASAPI32\ConsoleTracingMask
```



Reports as input file

Neurlux takes in input entire textual report, divided in tokens fed to deep neural networks / transformers

NEBULA: Self-attention for Dynamic Malware Analysis



Learning the language of malware
Joint work to understand how malware behaves from reports, leveraging cutting-edge transformers that outperforms other models in NLP (thank you Dmitrijs Trizna, national AI Ph.D. student & SSR @ Microsoft)

Static or dynamic?

Static > dynamic

Recent work show how much dynamic analysis is prone to error and noisy, opposed to the static one

... but dynamic better on unknowns

Static analysis perform worse on unseen families, while dynamic seems to be a bit better

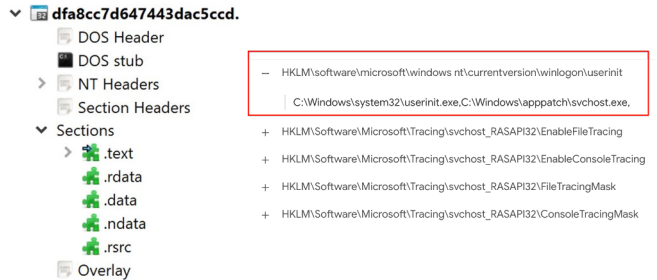
Static + dynamic?

The combination is not improving much the performances of static classifiers alone

Decoding the Secrets of Machine Learning in Windows Malware Classification: A Deep Dive into Datasets, Features, and Model Performance

Savino Dambra*	Yufei Han*	Simone Aonzo*	Platon Kotzias*
Norton Research Group	INRIA	EURECOM	Norton Research Group
Antonino Vitale	Juan Caballero	Davide Balzarotti	Leyla Bilge
EURECOM	IMDEA Software Institute	EURECOM	Norton Research Group

What about the best of ALL world?



Missing composite solution

Currently no proposal from the literature for complete AI systems that detect malware using rules, static, and dynamic analysis (also to mimic industry settings)

How to connect them?

While there are study that show that dynamic analysis improve accuracy by only few percentage points, there are still no clues on how to compose these layers

On-going work

We are investigating the performance of these AI systems to also improve the quality of testing (Thanks to Andrea Ponte again!)

Take-home messages of Part 1

Machine learning helps in the fight

Learning from data generalizes better than collecting rules and hashes

Domain knowledge vs end-to-end

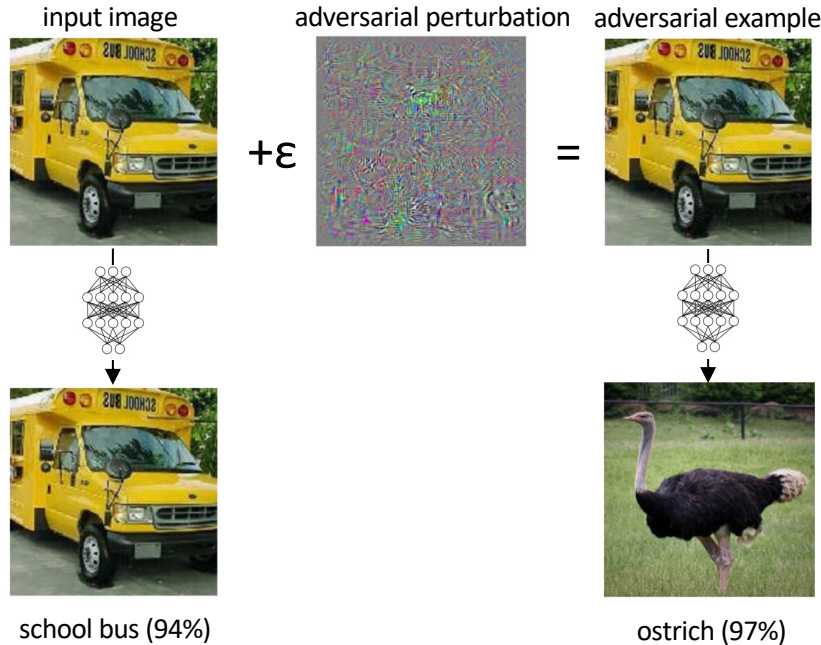
Data can be used as-is or with feature extraction, and the second performs better

Static vs Dynamic

Different ways to recognize malicious patterns, either from the structure of the file or from its behavior. Recent study shows that static is better.

Part 2 - Adversarial EXEmples

Adversarial Examples (Gradient-based Evasion Attacks)

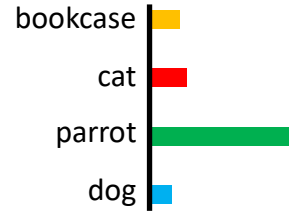
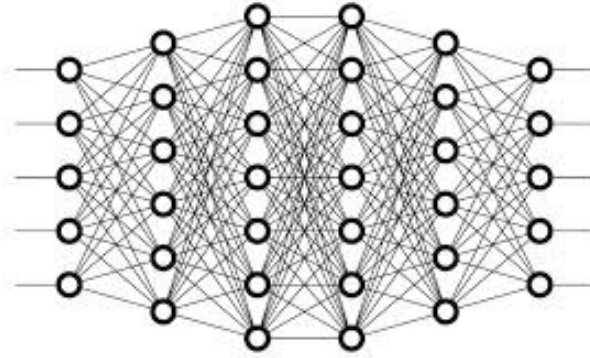


(Thanks to Battista Biggio, Maura Pintor for these slides!)

Szegedy et al., [Intriguing properties of neural networks](#), ICLR 2014

Biggio, Roli, et al., [Evasion attacks against machine learning at test time](#), ECML-PKDD 2013

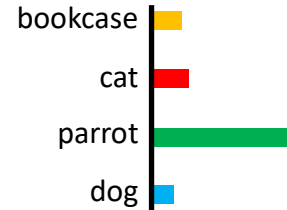
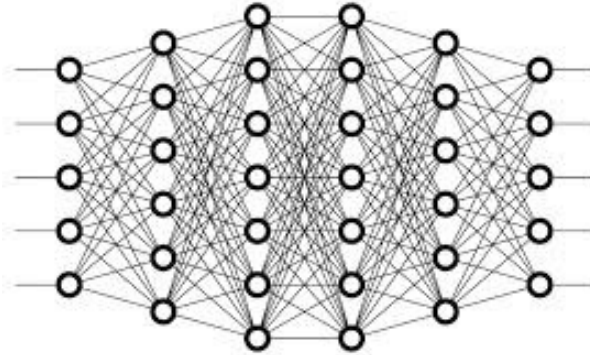
Adversarial Attacks



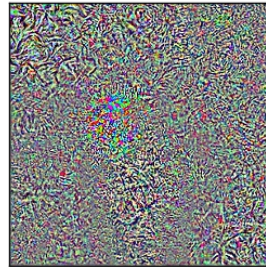
Adversarial attacks exploit the same underlying mechanism of learning, but aim to maximize the probability of error on the input data: $\max_D L(D; \mathbf{w})$

This problem can also be solved with gradient-based optimizers
(Biggio, Roli et al., ICML 2012; Biggio, Roli et al., ECML 2013; Szegedy et al., ICLR 2014)

How Do These Attacks Work?

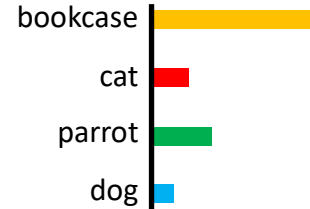
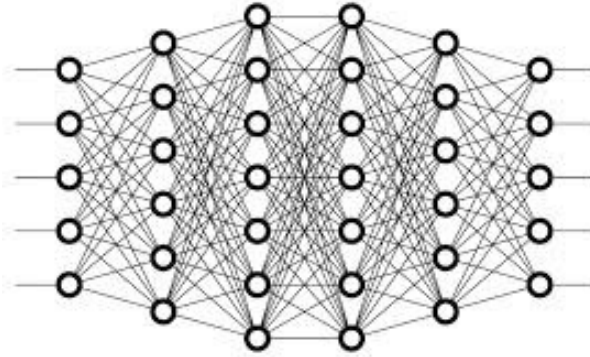


$$\max_D L(D; \mathbf{w})$$

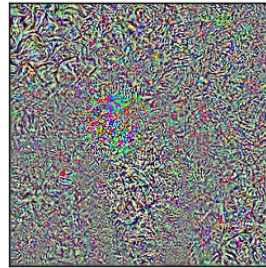


The gradient of the objective allows us to compute an *adversarial perturbation*...

How Do These Attacks Work?



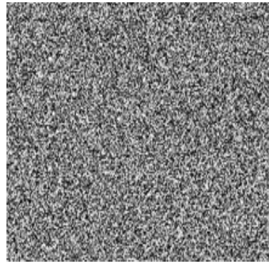
... which is then added to the input image to cause misclassification



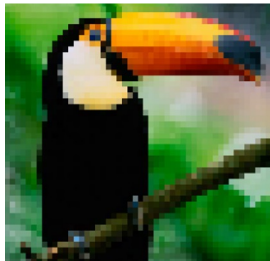
Same for malware?



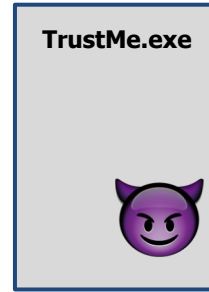
toucan (97%)



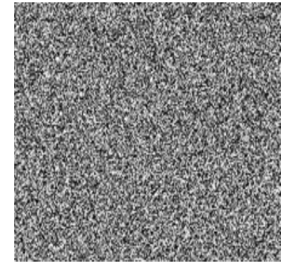
adversarial noise



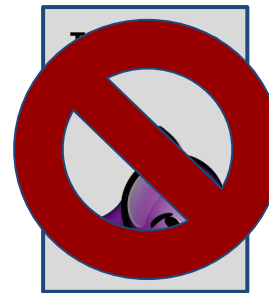
cat (95%)



malware (98%)



adversarial noise



Not runnable anymore!

Byte-sequences are not numbers

Programs and images are encoded in bytes

RGB is “continuous”, code instructions are not!

Distance between programs is undefined

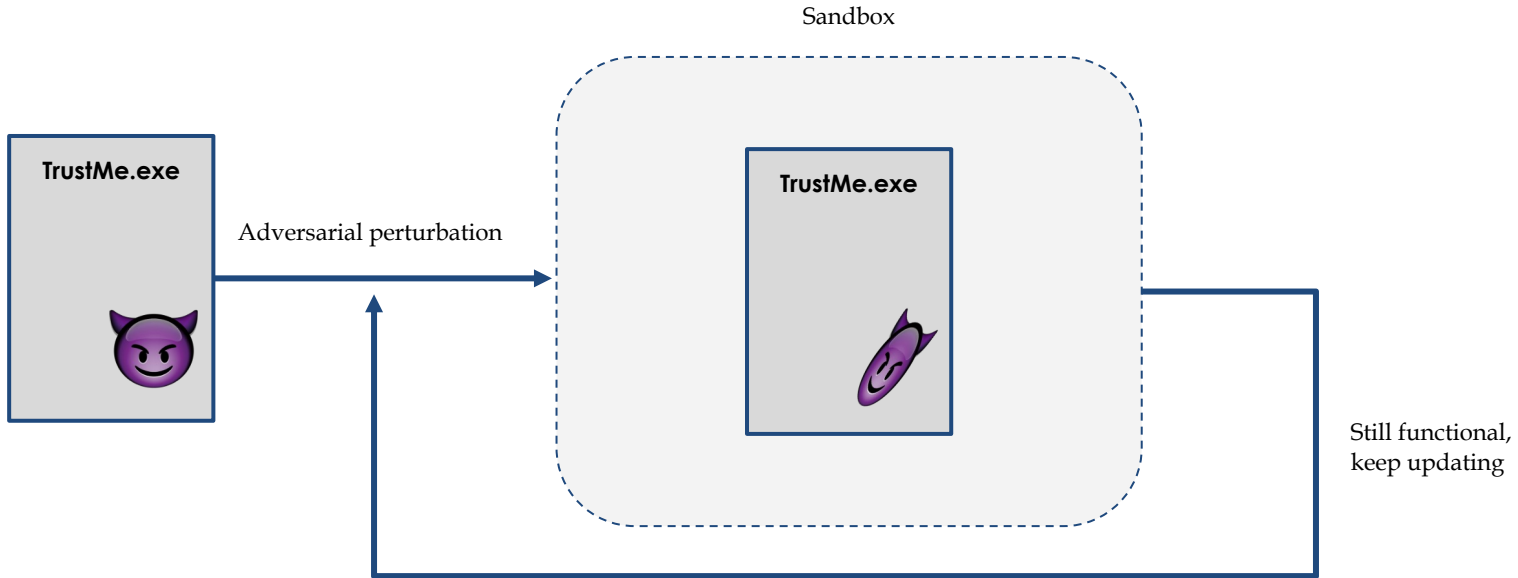
Example: ASCII table

What does $\sqrt{'a' - 'b'}$ means?

Dec	Hx	Oct	Char	Dec	Hx	Oct	Htmi	Chr	Dec	Hx	Oct	Htmi	Chr	Dec	Hx	Oct	Htmi	Chr
0	0	000	NUL (null)	32	20	040	€#32;	Space	64	40	100	€#64;	@	96	60	140	€#96;	`
1	1	001	SOH (start of heading)	33	21	041	€#33;	!	65	41	101	€#65;	A	97	61	141	€#97;	a
2	2	002	STX (start of text)	34	22	042	€#34;	"	66	42	102	€#66;	B	98	62	142	€#98;	b
3	3	003	ETX (end of text)	35	23	043	€#35;	#	67	43	103	€#67;	C	99	63	143	€#99;	c
4	4	004	EOF (end of transmission)	36	24	044	€#36;	\$	68	44	104	€#68;	D	100	64	144	€#100;	d
5	5	005	ENQ (enquiry)	37	25	045	€#37;	%	69	45	105	€#69;	E	101	65	145	€#101;	e
6	6	006	ACK (acknowledge)	38	26	046	€#38;	&	70	46	106	€#70;	F	102	66	146	€#102;	f
7	7	007	BEL (bell)	39	27	047	€#39;	'	71	47	107	€#71;	G	103	67	147	€#103;	g
8	8	010	BS (backspace)	40	28	050	€#40;	(72	48	110	€#72;	H	104	68	150	€#104;	h
9	9	011	TAB (horizontal tab)	41	29	051	€#41;)	73	49	111	€#73;	I	105	69	151	€#105;	i
10	A	012	LF (NL line feed, new line)	42	2A	052	€#42;	*	74	4A	112	€#74;	J	106	6A	152	€#106;	j
11	B	013	VT (vertical tab)	43	2B	053	€#43;	+	75	4B	113	€#75;	K	107	6B	153	€#107;	k
12	C	014	FF (NP form feed, new page)	44	2C	054	€#44;	,	76	4C	114	€#76;	L	108	6C	154	€#108;	l
13	D	015	CR (carriage return)	45	2D	055	€#45;	-	77	4D	115	€#77;	M	109	6D	155	€#109;	m
14	E	016	SO (shift out)	46	2E	056	€#46;	.	78	4E	116	€#78;	N	110	6E	156	€#110;	n
15	F	017	SI (shift in)	47	2F	057	€#47;	/	79	4F	117	€#79;	O	111	6F	157	€#111;	o
16	10	020	DLE (data link escape)	48	30	060	€#48;	0	80	50	120	€#80;	P	112	70	160	€#112;	p
17	11	021	DC1 (device control 1)	49	31	061	€#49;	1	81	51	121	€#81;	Q	113	71	161	€#113;	q
18	12	022	DC2 (device control 2)	50	32	062	€#50;	2	82	52	122	€#82;	R	114	72	162	€#114;	r
19	13	023	DC3 (device control 3)	51	33	063	€#51;	3	83	53	123	€#83;	S	115	73	163	€#115;	s
20	14	024	DC4 (device control 4)	52	34	064	€#52;	4	84	54	124	€#84;	T	116	74	164	€#116;	t
21	15	025	NAK (negative acknowledge)	53	35	065	€#53;	5	85	55	125	€#85;	U	117	75	165	€#117;	u
22	16	026	SYN (synchronous idle)	54	36	066	€#54;	6	86	56	126	€#86;	V	118	76	166	€#118;	v
23	17	027	ETB (end of trans. block)	55	37	067	€#55;	7	87	57	127	€#87;	W	119	77	167	€#119;	w
24	18	030	CAN (cancel)	56	38	070	€#56;	8	88	58	130	€#88;	X	120	78	170	€#120;	x
25	19	031	EM (end of medium)	57	39	071	€#57;	9	89	59	131	€#89;	Y	121	79	171	€#121;	y
26	1A	032	SUB (substitute)	58	3A	072	€#58;	:	90	5A	132	€#90;	Z	122	7A	172	€#122;	z
27	1B	033	ESC (escape)	59	3B	073	€#59;	;	91	5B	133	€#91;	[123	7B	173	€#123;	{
28	1C	034	FS (file separator)	60	3C	074	€#60;	<	92	5C	134	€#92;	\	124	7C	174	€#124;	
29	1D	035	GS (group separator)	61	3D	075	€#61;	=	93	5D	135	€#93;]	125	7D	175	€#125;	}
30	1E	036	RS (record separator)	62	3E	076	€#62;	>	94	5E	136	€#94;	^	126	7E	176	€#126;	~
31	1F	037	US (unit separator)	63	3F	077	€#63;	?	95	5F	137	€#95;	_	127	7F	177	€#127;	DEL

Source: www.asciitable.com

Preserve the original functionality

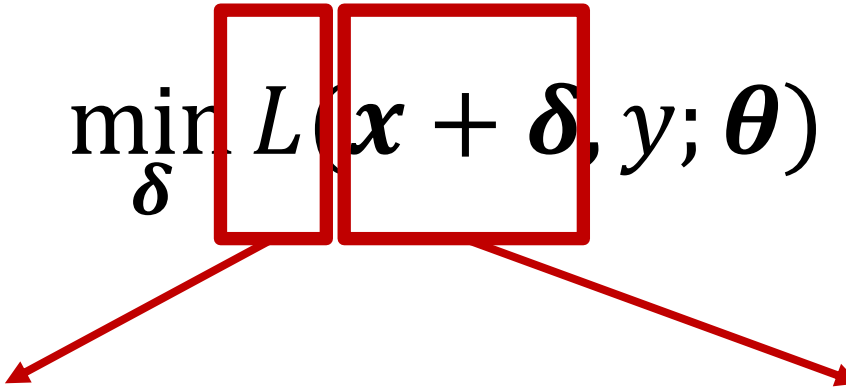


How to bridge these gaps?

1. Formulate the minimization problem differently
2. Study the format that represent programs
3. Understand how to exploit the format
4. Chose how to inject or perturb the content

Formulation of the problem

Adversarial attacks for images

$$\min_{\delta} L(\mathbf{x} + \delta, y; \theta)$$


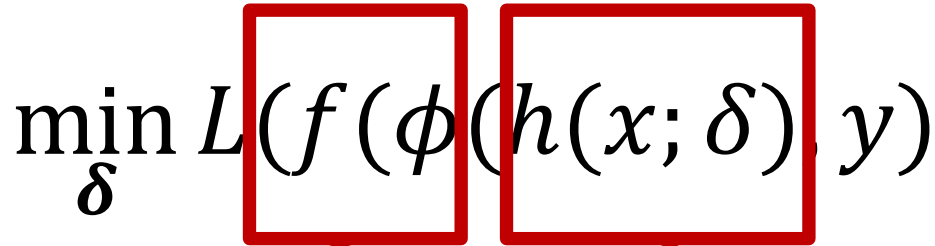
Network architecture in the loss

All the internals of a neural network / shallow model are hidden inside the loss

Additive Manipulation

Input samples are injected with additive noise, without any concern on the structure of the file

Adversarial attacks for security detectors

$$\min_{\delta} L(f(\phi(h(x; \delta)), y))$$


Model function and features

Need to explicit the model function and the features, since they might be non differentiable

Practical Manipulations

No additions, but a complex function that handles format specification by design

Take-home message: implementing an attack

$$\min_{\delta} L(f(\phi(h(x; \delta)), y))$$

Define the Optimizer
Depending on the differentiability of the components, pick a gradient-based or gradient-free algorithm

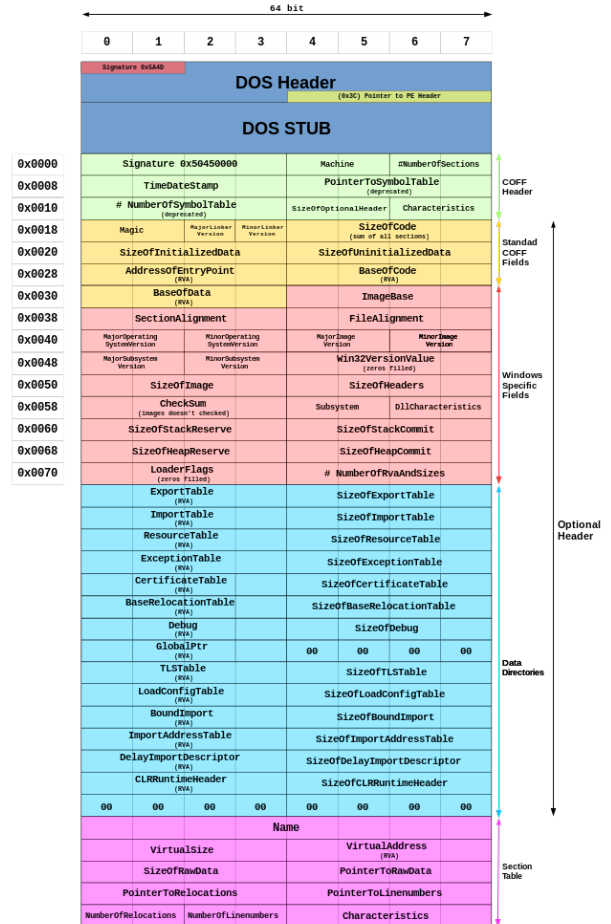
Define the Manipulations
Study the format, understand its ambiguities, and write manipulations that do not break the original functionality

Manipulations of Windows PE file format

Windows PE File Format

Format adapted for “modern” programs
(from Windows NT 3.1 on)

Before there were other formats, one is the
DOS (kept for retrocompatibility)

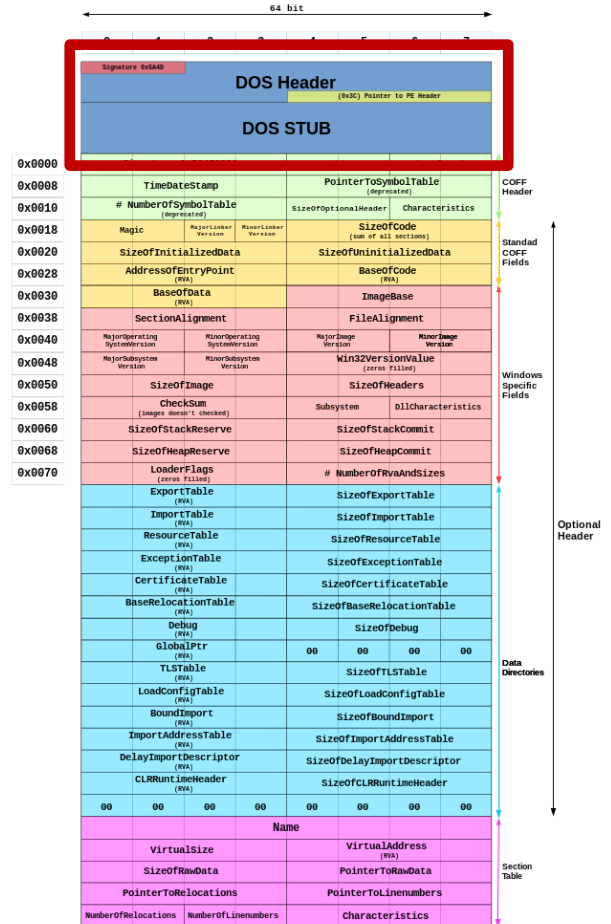


Windows PE File Format

DOS Header + Stub

Metadata for DOS program

Executing a modern program in DOS will trigger the “This program cannot be run in DOS mode” output

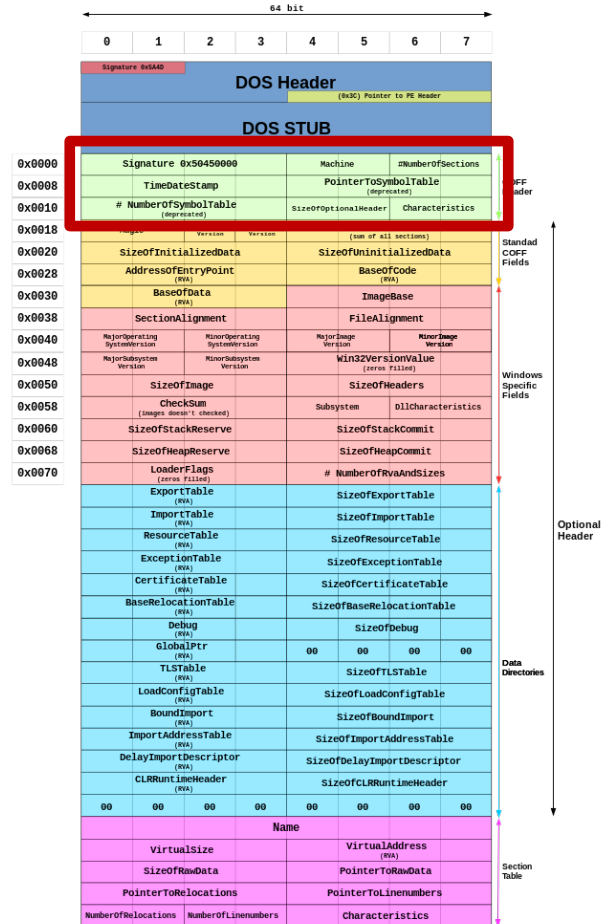


Windows PE File Format

PE Header

Real metadata of the program

Describes general information of the file

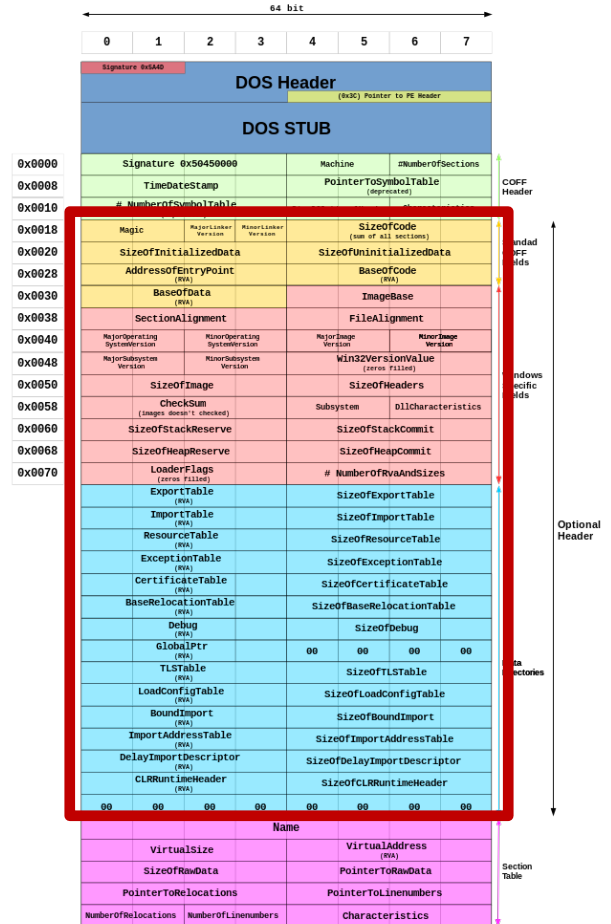


Windows PE File Format

Optional Header

Spoiler: not optional at all :)

Instructs the loader where to find each object inside the file



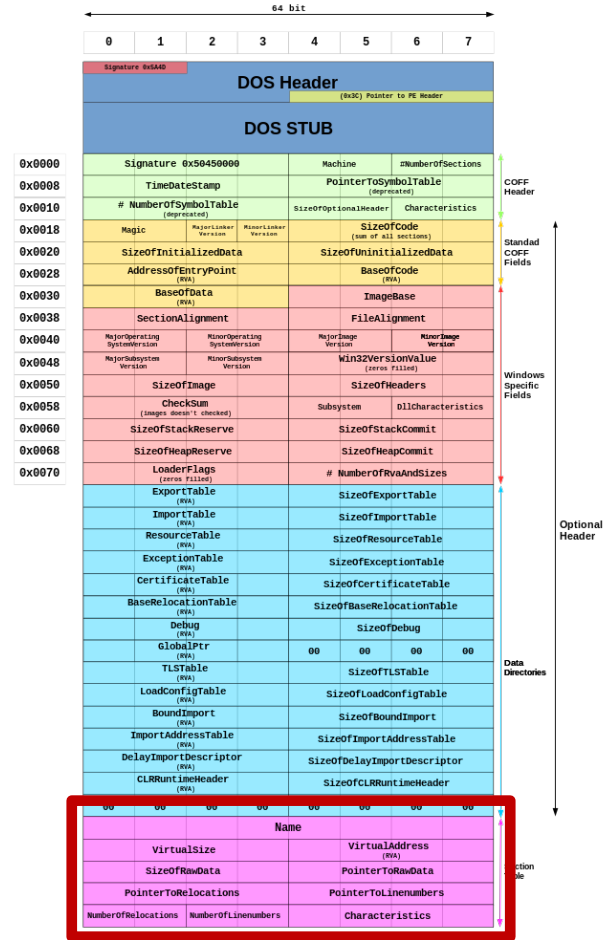
Windows PE File Format

Section Table and Sections

Describes where to find code, initialized data, resources, etc to the loader

These are “sections”, and each has a “section entry” with its characteristics

Examples: code is “.text”, read-only data is “.rodata”, resources are “.rsc”, and counting



How programs are loaded

1 Headers

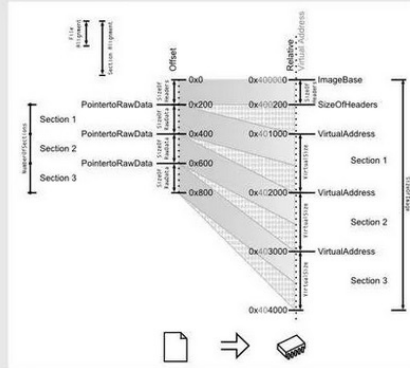
the *DOS Header* is parsed
the *PE Header* is parsed
(its offset is *DOS Header's e_lfanew*)
the *Optional Header* is parsed
(it follows the *PE Header*)

2 Sections table

Sections table is parsed
(it is located at: $\text{offset}(\text{OptionalHeader}) + \text{SizeOfOptionalHeader}$)
it contains *NumberOfSections* elements
it is checked for validity with alignments:
FileAlignments and *SectionAlignments*

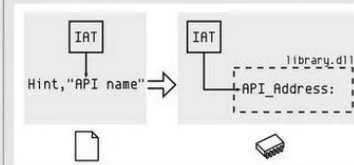
3 Mapping

the file is mapped in memory according to:
the *ImageBase*
the *SizeOfHeaders*
the *Sections* table



4 Imports

DataDirectories are parsed
they follow the *OptionalHeader*
their number is *NumOfRVAAndSizes*
imports are always #2
Imports are parsed
each descriptor specifies a *DLLname*
this DLL is loaded in memory
IAT and *INT* are parsed simultaneously
for each API in *INT*
its address is written in the *IAT* entry



5 Execution

Code is called at the *EntryPoint*
the calls of the code go via the *IAT* to the APIs



<https://code.google.com/archive/p/corkami/wikis/PE101.wiki>

Towards Adversarial EXEmples

Perturb the representation of a file

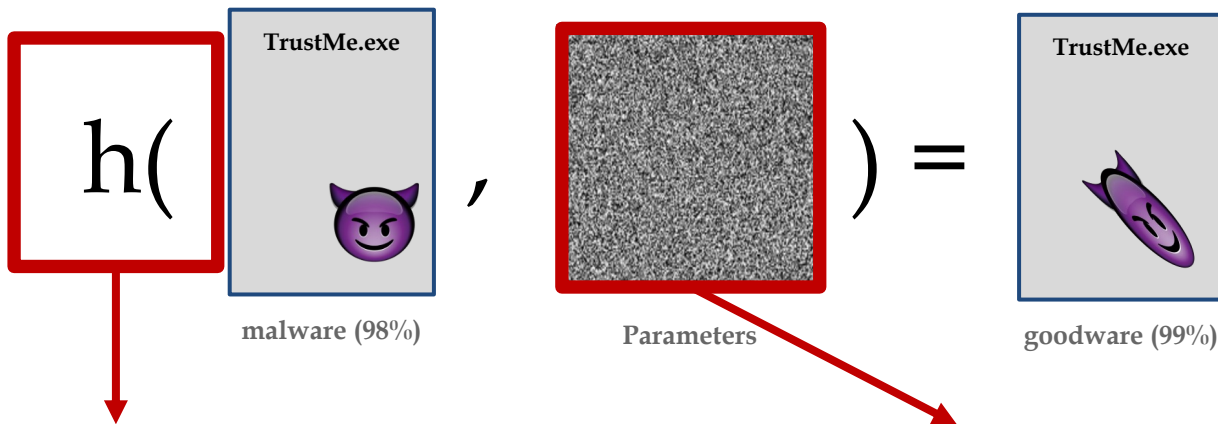
Keep intact the original functionality

Example: rotation for images

How to bridge the gap?



Practical Manipulations



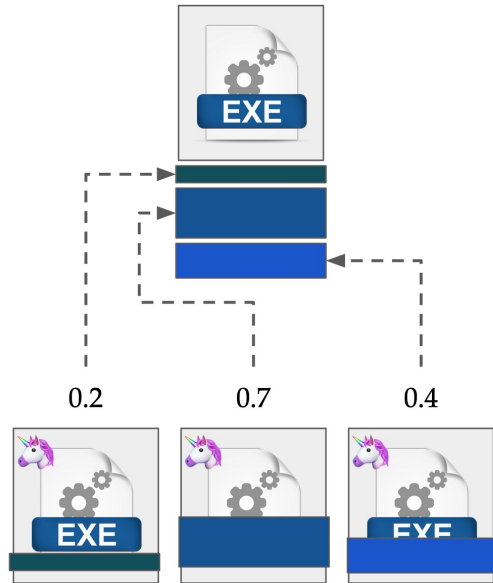
Practical Manipulation function

Alter file representation
without destroying the structure
and the functionalities and avoid
usage of sandboxes

Parameters

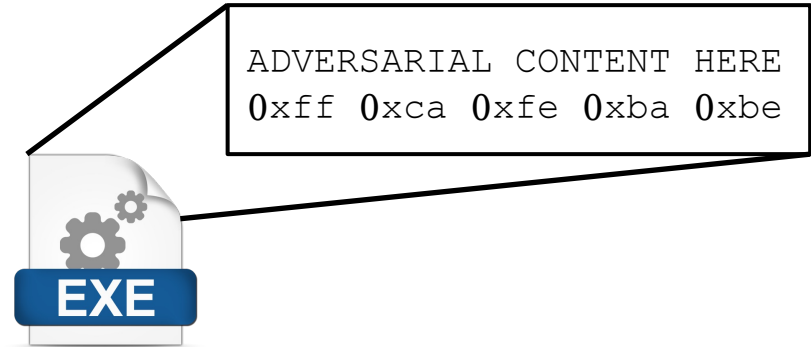
Manipulations are parametrized
so an optimization algorithm can
fine tune them

Structural Manipulations



Injecting content

Alter file structure to include more byte sequences



Replacing content

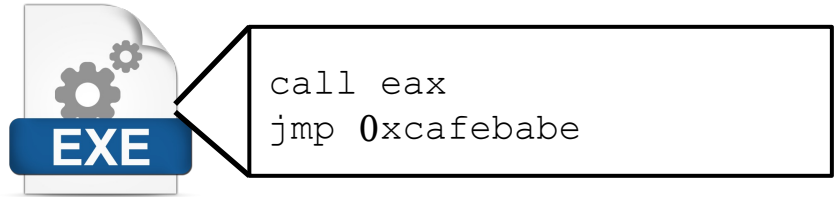
Leverage ambiguous file format specifications to alter bytes that are not considered at runtime

Behavioral Manipulations

WARNING
more difficult to implement!



Packing and obfuscation
Encrypt program inside another one, or complicate the sequence of instructions



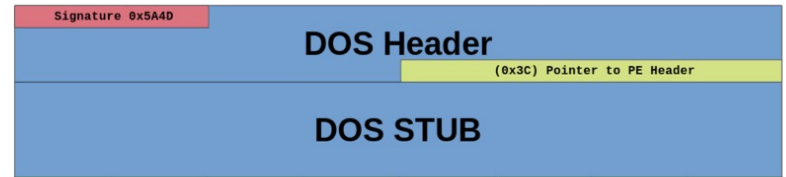
Inject new execution flows
Call APIs, add loops, jump to new code sections, and more

DOS header perturbations

The attacker edit as many bytes as they want

Untouched: magic number MZ and offset to real PE header

Content loaded in memory, not executed

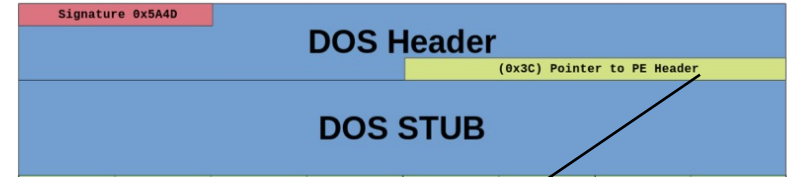


DOS header extension

Exploit offset to real header, increment value

Insert arbitrary content between DOS header and PE header

Content loaded into memory, not executed



Signature 0x50450000	Machine	#NumberOfSections
TimeDateStamp	PointerToSymbolTable (deprecated)	
# NumberOfSymbolTable (deprecated)	SizeOfOptionalHeader	Characteristics

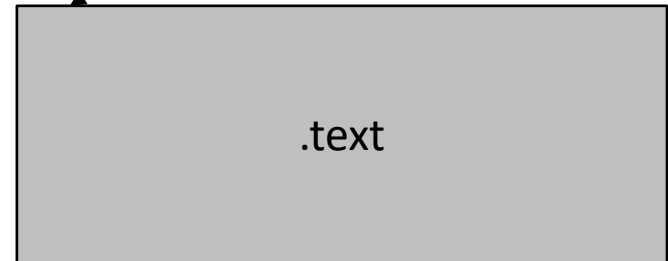
Content shifting

Exploit offset in section entry, increment to manipulate the loader in searching for section content

The attacker can inject content after the section table, or between sections

NOT LOADED IN MEMORY, skipped by the loader

		Name			
VirtualSize				VirtualAddress (RVA)	
SizeOfRawData				PointerToRawData	
PointerToRelocations				PointerToLinenumbers	
NumberOfRelocations	NumberOfLinenumbers			Characteristics	



Section Injection

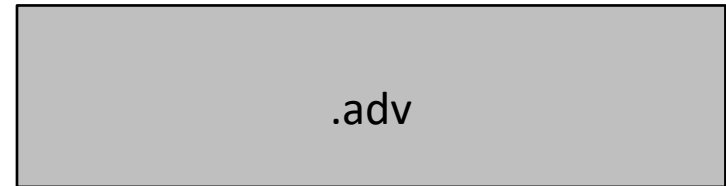
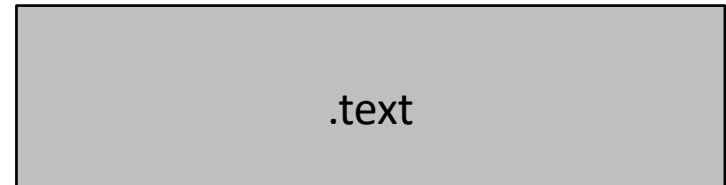
Manipulate section table to add new entry

Append chunk of bytes, referenced by newly added entry

Loaded in memory or not, depending by the characteristics set up inside the entry

		Name		
VirtualSize			VirtualAddress	(RVA)
SizeOfRawData			PointerToRawData	
PointerToRelocations			PointerToLinenumbers	
NumberOfRelocations	NumberOfLinenumbers		Characteristics	

		Name		
VirtualSize			VirtualAddress	(RVA)
SizeOfRawData			PointerToRawData	
PointerToRelocations			PointerToLinenumbers	
NumberOfRelocations	NumberOfLinenumbers		Characteristics	

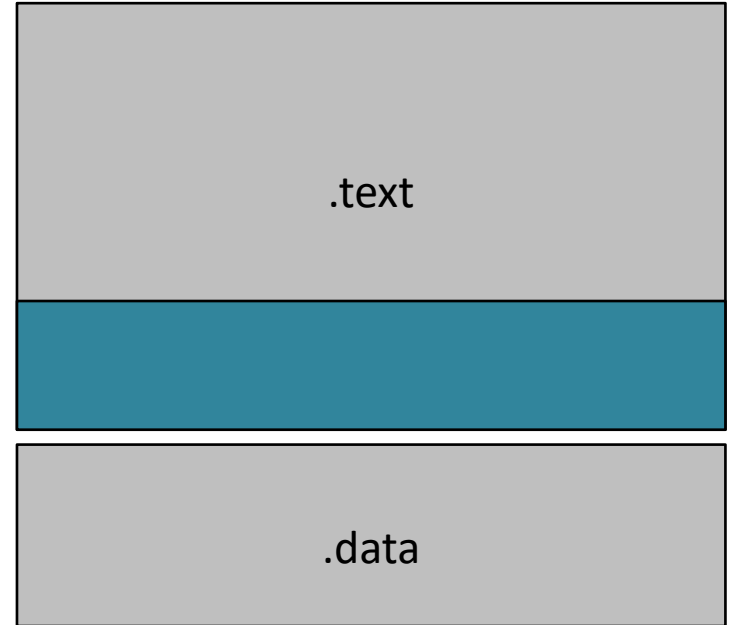


Slack space

Section content is padded with 0 to keep file alignments

The attacker can rewrite such slack space

Loaded in memory, not executed



Padding

Appending content at the end

Most trivial manipulation

Not loaded in memory



Optimization Algorithms

Choosing the strategy accordingly

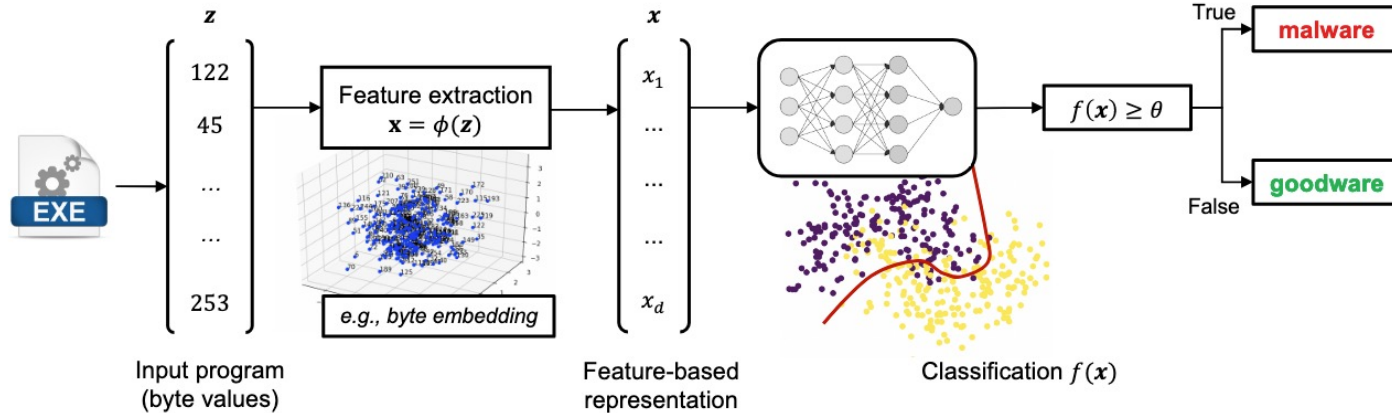
Gradient-based

Own the model
AND
Model is differentiable

Gradient-free

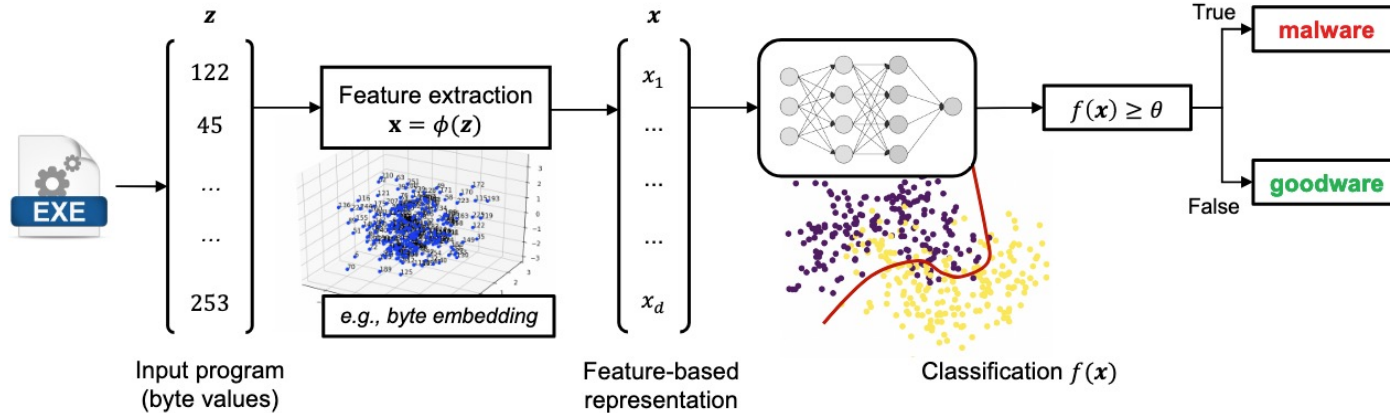
Model not accessible
OR
Model is not differentiable

Gradient-based strategies



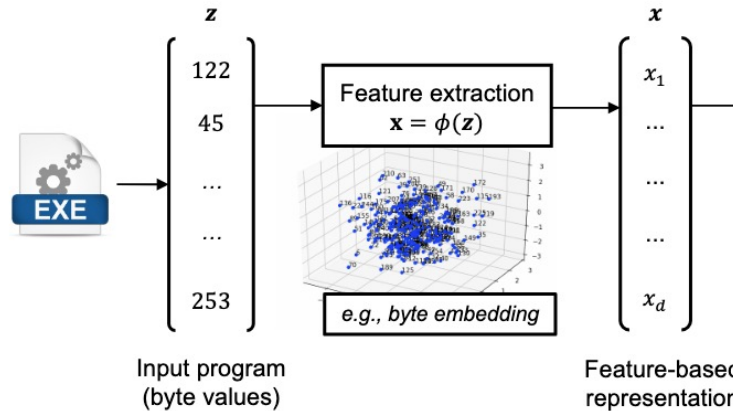
Use gradient-descent to compute adversarial examples (as for images)

Gradient-based strategies



~~Use gradient descent to compute adversarial examples (as for images)~~
Bytes do not have a distance metric, a feature extractor is **ALWAYS** needed to compute something meaningful

Embedding for end-to-end networks



All bytes are replaced with a vector learned at training time,
where a distance metric is imposed...

... but the embedding layer is **not differentiable**

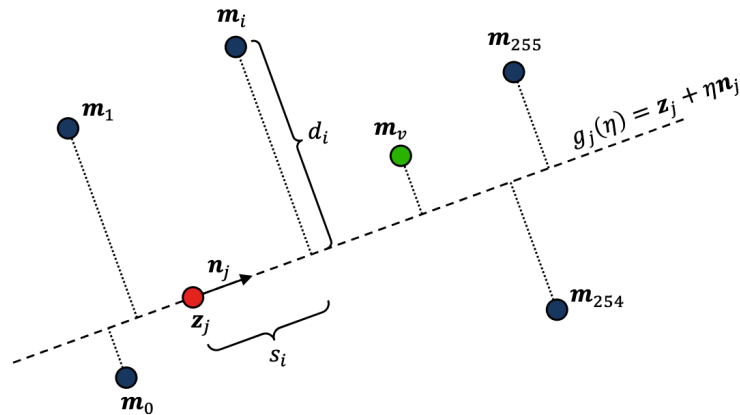
How to propagate gradient information?

$$\frac{\partial L}{\partial \delta} = \frac{\partial L}{\partial f} \frac{\partial f}{\partial \phi} \frac{\partial \phi}{\partial h} \frac{\partial h}{\partial \delta}$$

End-to-end gradient
you would like to
compute

**Non-differentiable
manipulations and
embedding!**

Solution: change the optimizer



Still gradient descent, but inside the **embedding space!**

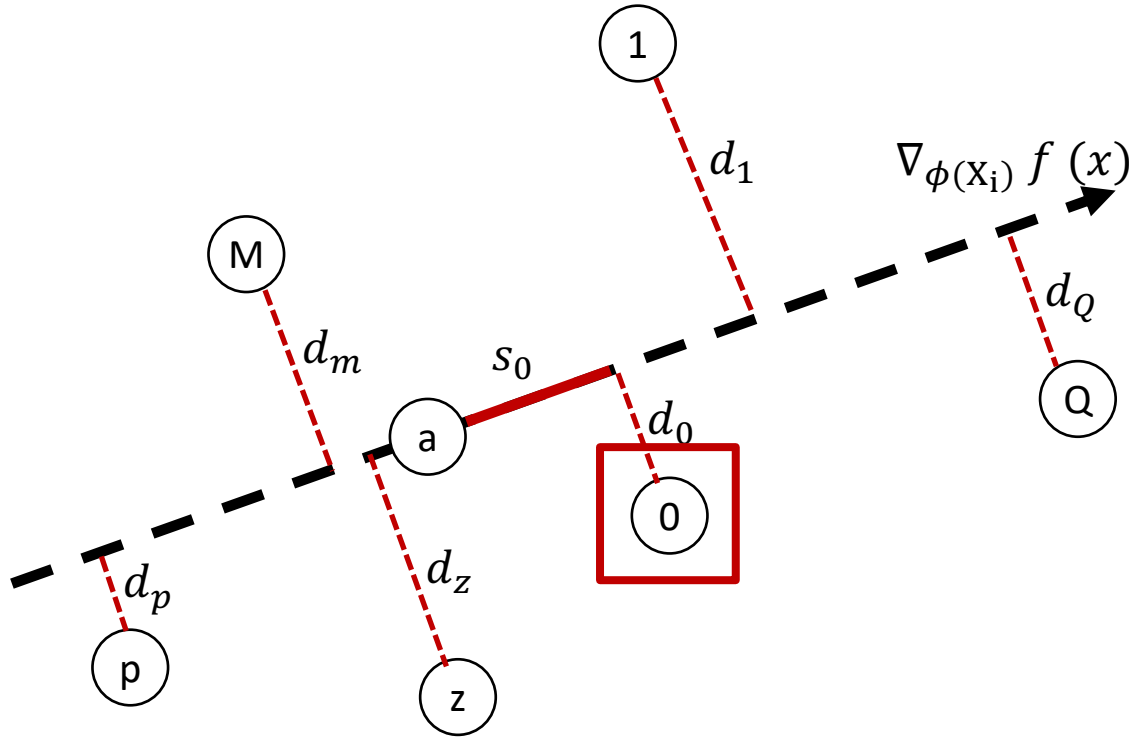
Optimize where gradients are available and reconstruct bytes **after** the search

BGD: Byte Gradient Descent

```
1  $\mathbf{E}_i = \hat{\phi}(i), \forall i \in [0, 256]$ 
2  $\mathbf{t}^{(0)} \in \mathcal{T}$ 
3 for  $i$  in  $[0, N - 1]$ 
4    $\mathbf{X} \leftarrow \phi(h(\mathbf{z}, \mathbf{t}^{(i)}))$  // feat. space
5    $\mathbf{G} \leftarrow -\nabla_{\mathbf{X}} f(\mathbf{X}) \odot \mathbf{m}$ 
6    $\mathbf{g} \leftarrow (\|\mathbf{G}_0\|, \dots, \|\mathbf{G}_n\|)$ 
7   for  $k$  in  $\text{argsort}(\mathbf{g})_{0, \dots, \gamma} \wedge g_k \neq 0$ 
8     for  $j$  in  $[0, \dots, 255]$ 
9        $\mathbf{S}_{k,j} \leftarrow \mathbf{G}_k^t \cdot (\mathbf{E}_j - \mathbf{X}_k)$ 
10       $\tilde{\mathbf{X}}_{k,j} \leftarrow \|\mathbf{E}_j - (\mathbf{X}_k + \mathbf{G}_k \mathbf{S}_{k,j})\|_2$ 
11       $\mathbf{t}_k^{(i+1)} \leftarrow \arg \min_{j: \mathbf{S}_{k,j} > 0} \tilde{\mathbf{X}}_{k,j}$  //input space
12  $\mathbf{t}^{\star} \leftarrow \mathbf{t}^{(N)}$ 
13  $\mathbf{z}^{\star} \leftarrow h(\mathbf{z}, \mathbf{t}^{\star})$ 
14 return  $\mathbf{z}^{\star}$ 
```

1. Compute gradient in feature space
2. Define a way for replacing values
For bytes: inverse look-up of embedding
3. Follow the direction of gradient and replace byte with other byte

BGD: Byte Gradient Descent (optimization)



The process is repeated according to the **stepsize of the attack**, that quantifies how many bytes are modified at each iteration

BGD: Byte Gradient Descent (reconstruction)

At the end of each iteration, I need to replace one byte, **not an embedding value**

But each byte is chosen in the embedding space, reconstruction just invert the look-up function

$$\arg \min_{j: S_{k,j} > 0} \tilde{X}_{k,j}$$

Choosing the strategy accordingly

Gradient-based

Own the model

AND

Model is differentiable

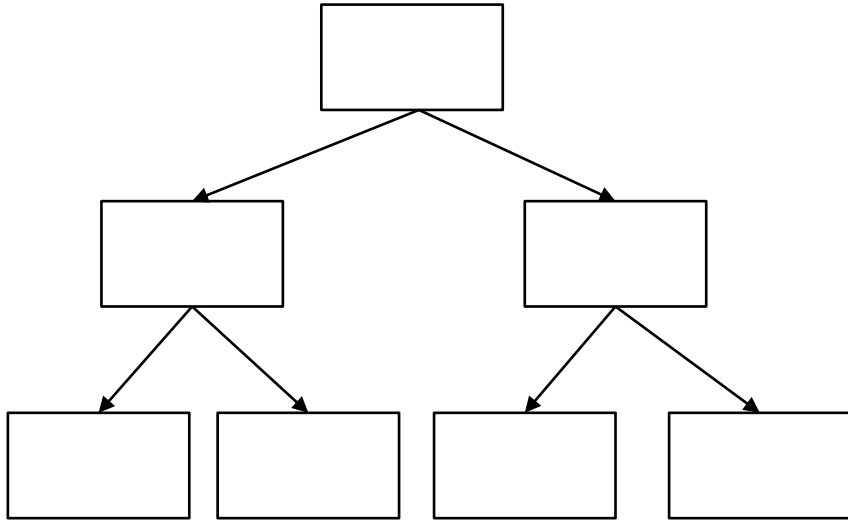
Gradient-free

Model not accessible

OR

Model is not differentiable

Reality check: robust models are not differentiable



State-of-the-art classifiers use decision trees

No gradients can be computed

Reality check: most models are unavailable

61 / 68 security vendors flagged this file as malicious

4c1dc737915d76b7ce579abddaba74ead6fbb5b519a1ea45308b8c49b950655c
4c1dc737915d76b7ce579abddaba74ead6fbb5b519a1ea45308b8c49b950655c.exe

788.00 KB Size | 2021-09-06 11:22:57 UTC | 3 days ago

EXE

Community Score

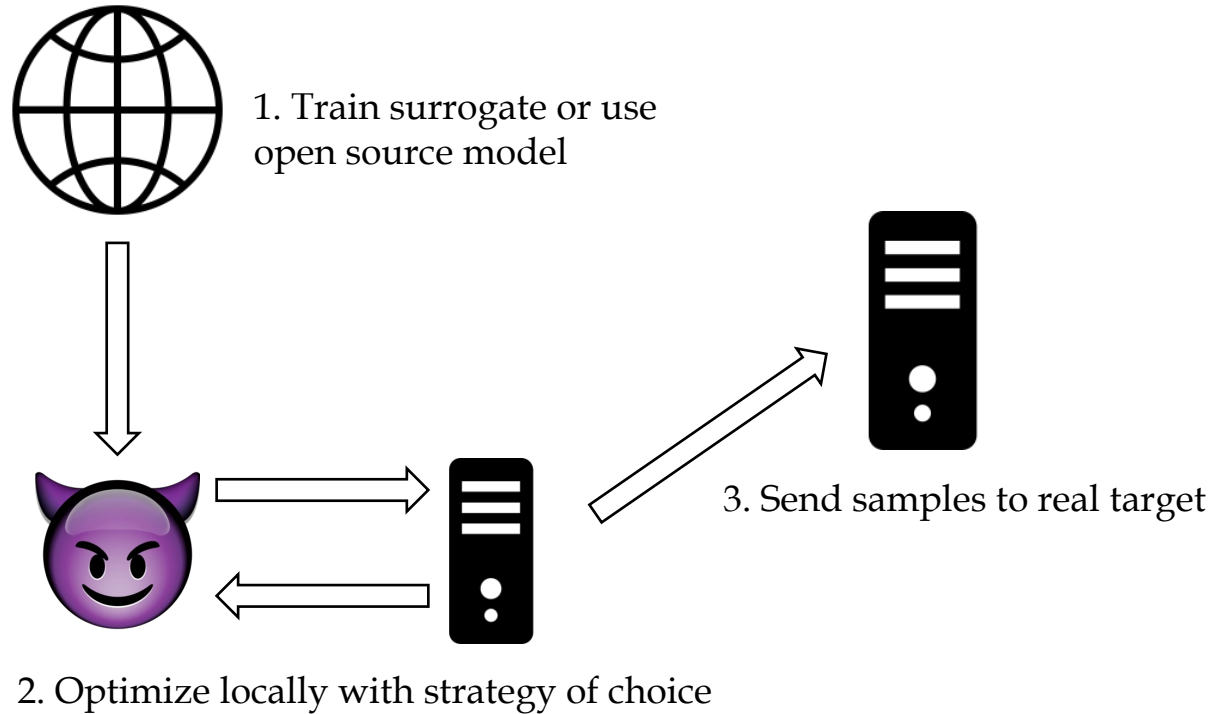
DETECTION	DETAILS	RELATIONS	BEHAVIOR	COMMUNITY
Ad-Aware	Trojan.Ransom.AUC		AhnLab-V3	Malware/Win32.RL_Generic.R295351
Alibaba	Ransom.Win32/Petya.404bad21		ALYac	Trojan.Ransom.Petya
SecureAge APEX	Malicious		Arcabit	Trojan.Ransom.AUC
Avast	Win32/Patched-AWP [Trj]		AVG	Win32/Patched-AWP [Trj]
Avira (no cloud)	TRAD.Petya.Y.hhcl		BitDefender	Trojan.Ransom.AUC
BitDefenderTheta	Gen:NN.Zexaf.34126.XuW@ay8Hrybt		Bkav Pro	W32.AIDetect.malware2
CAT-QuickHeal	Ransom.Petya.MUE.S6		ClamAV	Win.Trojan.Petya-6312160-0
Comodo	Malware@#304z9hhvmp31		CrowdStrike Falcon	Win/malicious_confidence_100% (W)
Cylance	Unsafe		Cynet	Malicious (score: 100)
Cyren	W32/Trojan.XMFF-8835		DrWeb	Trojan.MBRlock.245
eGambit	Unsafe.AI_Score_99%		Elastic	Malicious (high Confidence)

Most models are hosted on private servers

Detection performed in cloud

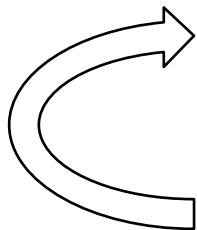
No gradients can be computed

Transfer attacks

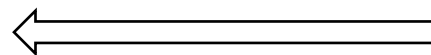
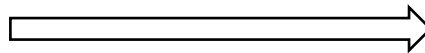


Query attacks

3. Perturb bytes of the sample, considering the scores from remote



1. Send sample to target



2. Obtain scores from remote



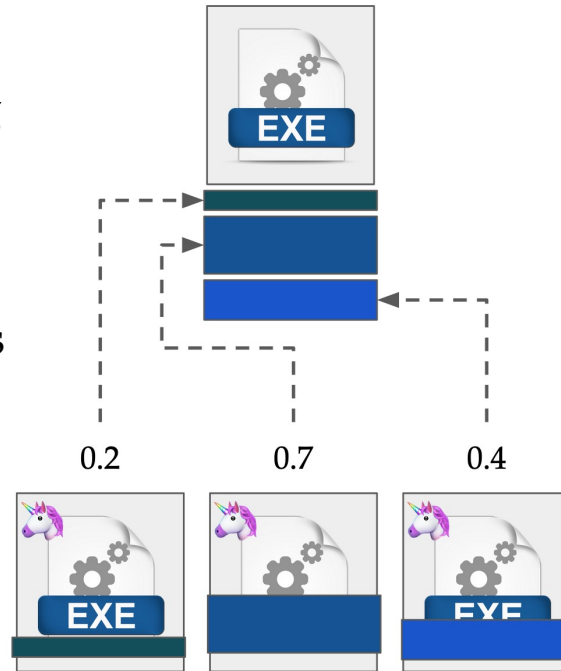
Very slow if optimizer works byte-per-byte

GAMMA: Speeding up by injecting benign content

Intuition

classifiers can be fooled by introducing content of the goodware class!

The optimizer explore less space, no modification byte-per-byte, but it relies on portions of goodware programs injected with practical manipulations



GAMMA against commercial products

	Malware	Random	Sect. Injection
AV1	93.5%	85.5%	30.5%
AV2	85.0%	78.0%	68.0%
AV3	85.0%	46.0%	43.5%
AV4	84.0%	83.5%	63.0%
AV5	83.5%	79.0%	73.0%
AV6	83.5%	82.5%	69.5%
AV7	83.5%	54.5%	52.5%
AV8	76.5%	71.5%	60.5%
AV9	67.0%	54.5%	16.5%

Breaking signatures and patterns
GAMMA (transfer) reduces the
performance of commercial products
hosted on VirusTotal!

Sneak preview 🕶️ (1/2)

Genetic algorithms are slow

Optimizing EXEmples with GAMMA requires plenty of time, and it is not easy to control the injected content

Zero-order optimization joins the fight

We are currently working on bringing zero-order optimization inside the world of EXEmples, bending the theory in this non-sensical world without metrics (thank you Marco Rando, Ph.D. student @ MALGA)

$$g_{(G,h)}(x) = \frac{d}{\ell} \sum_{i=1}^{\ell} \frac{f(x + hGe_i) - f(x - hGe_i)}{2h} Ge_i.$$

From theoretical guarantees to EXEmples

First results suggest an improvement in gradient-free optimization attacks against MalConv and GBDT. Stay tuned for interesting results, and improved optimization algorithms!

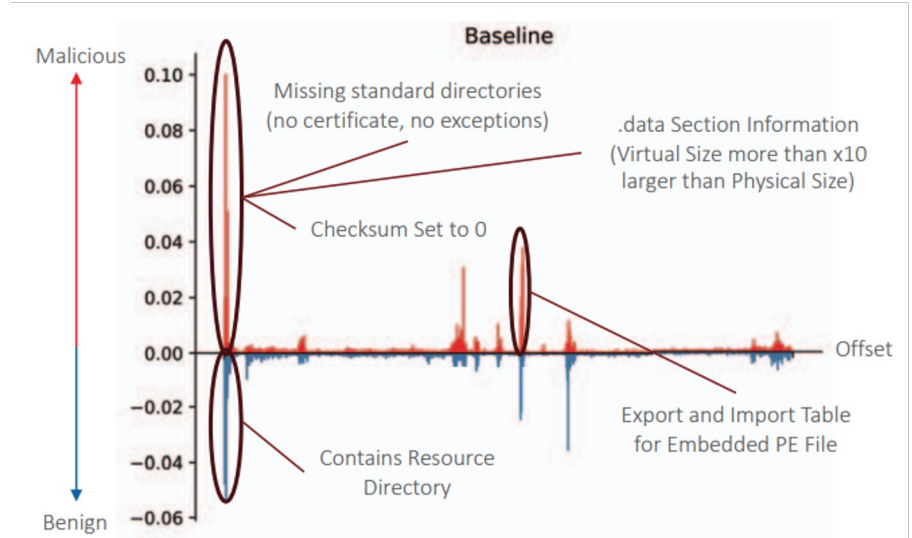
What are these models learning?

Great performances, but why?

Accuracy is high, false positives are low, but most machine learning models are difficult / impossible to inspect

Explainable AI

Train interpretable models (linear, trees) or apply explainability techniques to demystify decisions

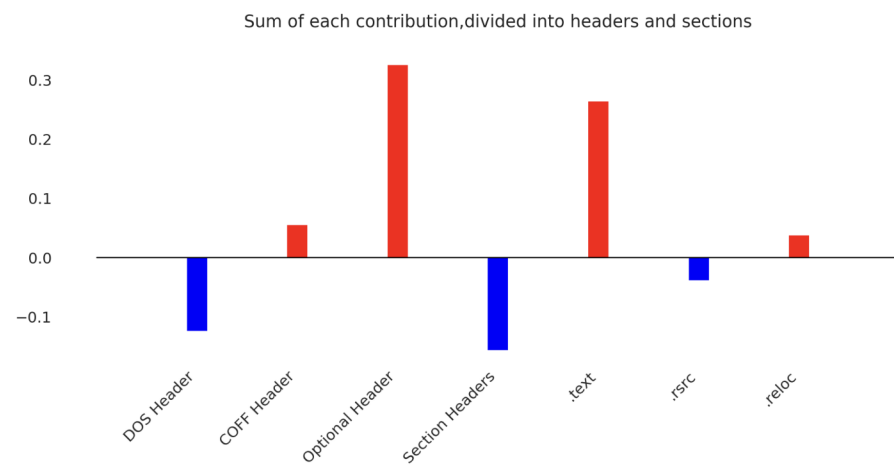


What are these models learning?

Sometimes, we don't know!

Our intuition on data is completely different from the correlation learnt by machine learning models

Example: malware detector attributes “legitimate” importance to unused space inside programs!



Sneak preview 🕶️ (2/2)

What about poisoning?

There are plenty of work on evasive EXEmples, but only few on poisoning of malware detectors

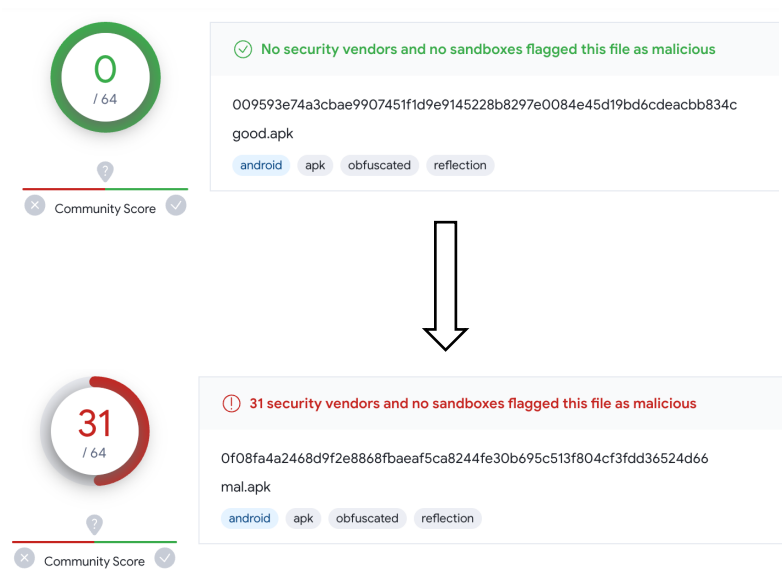
Your PAINT is now an EXEmple

Realistic scenario: attackers embed malicious signatures in regular software, since EVERYBODY trusts VirusTotal
Joint work to show the devastating effect of the exploitation of labelling systems

(thank you Simone Aonzo, Associate Professor @ EURECOM, Han Yufei, Senior Researcher @ INRIA, Tianwei Lan, Ph.D.)

Starting from Android, on EXE is easy

We are working on a more challenging scenario, which is embedding malicious objects inside Android applications
Stay tuned for interesting results!



Take-home messages of Part 2

New paradigm: study the format first

Not possible to re-use the same strategies, but first attackers must know how to deal with complex data structures

Adapt already-developed optimization algorithms

Not possible to re-use the same algorithms, since models might be only partially differentiable

Benign content injection rocks

Reduce the search space, faster attacks with effective results

Effectiveness in the real world as well

Evidence show that commercial products might be evaded as well

Are these model learning something?

Yes, but this is not what we expect, and spurious correlations are around the corner

Part 3: How to defend from EXEmples?

Recap: Adversarial EXEmples

Minimal byte perturbations

Many examples on how machine learning malware detectors can be bypassed with carefully-crafted input

Ambiguities of file format

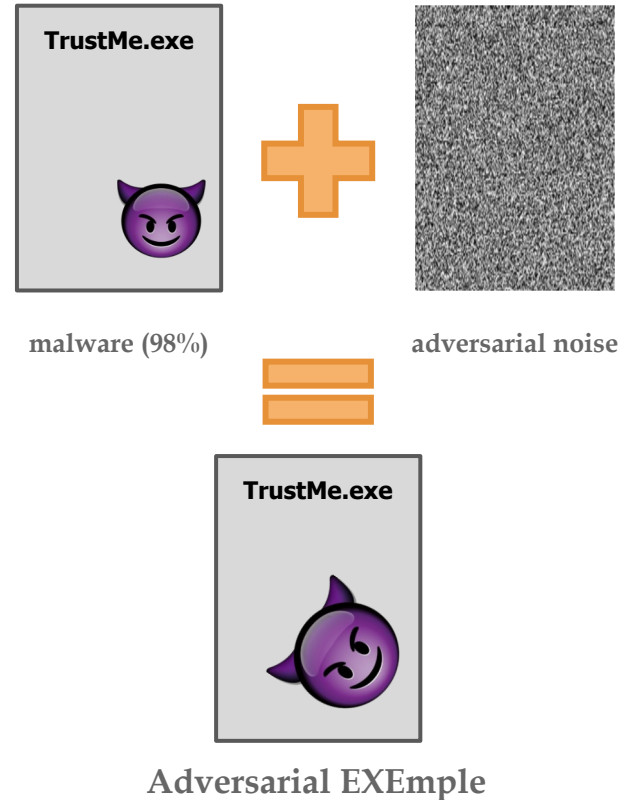
The Windows PE file format is redundant and many components are not used by the operating system at loading time, giving space to the attacker

Math is unreliable

The domain is discrete, models work mostly with continuous values, and attackers can “fill the blanks” with adversarial noise

How to avoid EXEmples?

Not clear how to patch this problem, but we isolated 4 relevant “claimed-to-be” robust malware detectors



Heuristic defense

Combination of pre-processing

Detect trivial manipulation, and then process input with ensemble of models

Partially reproducible

There are no pre-trained available, but code is available online

<https://github.com/EQuiw/2020-evasion-competition>

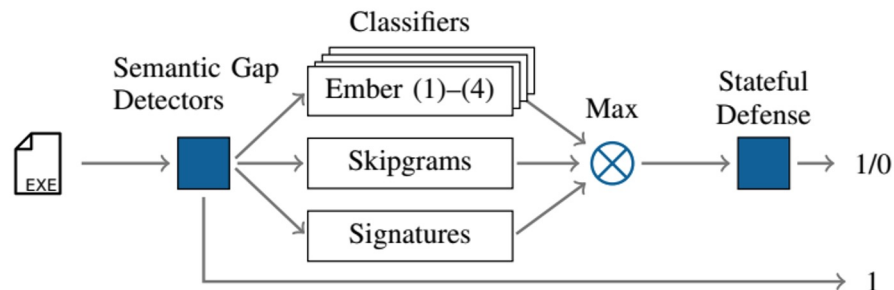
Unpublished

Still a preprint, never published (proposed for a competition)

Against All Odds: Winning the Defense Challenge in an Evasion Competition with Diversification

Erwin Quiring, Lukas Pirch, Michael Reimsbach, Daniel Arp, Konrad Rieck

Technische Universität Braunschweig
Braunschweig, Germany



Adversarial Training

$$\min_{\theta} \rho(\theta), \quad \text{where} \quad \rho(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{\delta \in \mathcal{S}} L(\theta, x + \delta, y) \right]$$

Train with EXamples

Computing state-of-the-art attacks and include them inside the training set (process is repeated until the achievement of the desired robustness)

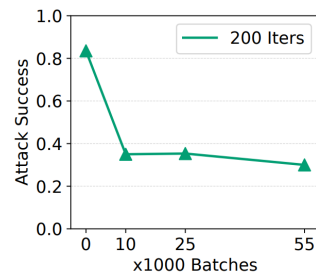
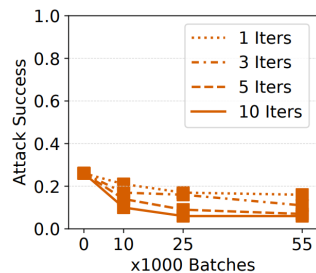
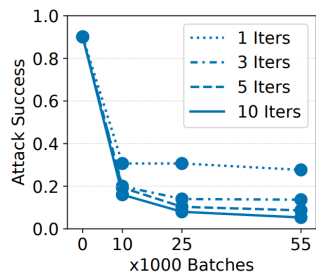
Adversarial Training

Alter code

Leverage behavioral manipulations to rewrite part of assembly code

Published, not reproducible

There are no pre-trained available, nor public source code that can be used to train a model. The technique is known, but the attacks used for this paper as well are closed



Adversarial Training for Raw-Binary Malware Classifiers

Keane Lucas
Carnegie Mellon University

Samruddhi Pai
Carnegie Mellon University

Weiran Lin
Carnegie Mellon University

Lujo Bauer
Carnegie Mellon University

Michael K. Reiter
Duke University

Mahmood Sharif
Tel Aviv University

Non-negative Networks

Malicious contributions

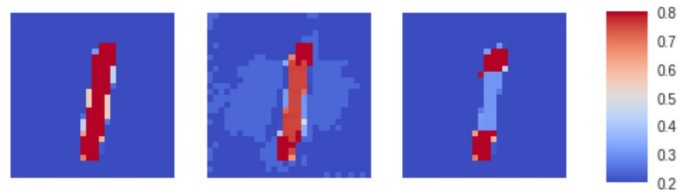
Intuition: classification is based on addition of small malicious triggers, until a threshold is reached

Remove negative weights

Train an end-to-end model, by clipping to positive values all the weights of the network

Attacks constrained

On images, non-negative network force attacks to only tamper with meaningful information



Non-negative Networks

Pre-trained available

Testing through model trained on EMBER,
released for a challenge
(performances are debatable)

Unpublished

Still a preprint, never published to either
conferences, journals or workshops

Hardly reproducible

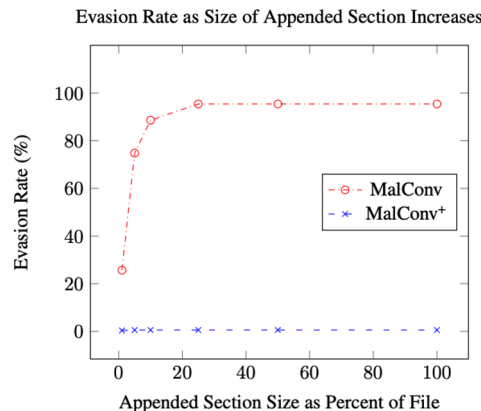
Even by trying to re-write code, many papers
tried and failed to train NonNeg MalConv

Non-Negative Networks Against Adversarial Attacks

William Fleshman,¹ Edward Raff,^{1,2} Jared Sylvester,^{1,2} Steven Forsyth,³ Mark McLean¹

¹Laboratory for Physical Sciences, ²Booz Allen Hamilton, ³Nvidia

{william.fleshman, edraff, jared, mrmclea}@lps.umd.edu, sforsyth@nvidia.com



Monotonic Classifiers

IWSPA'18, March 21, 2018, Tempe, AZ, USA

Malicious contributions

Intuition: classification is based on addition of small malicious triggers, until a threshold is reached

Gradient boosting decision tree

Use custom training process that trains an additive decision function

Subset of features

Not using EMBER, but the authors propose a reduced features set that is harder to manipulate

Not reproducible

There are no pre-trained available, nor public source code that can be used to train a model

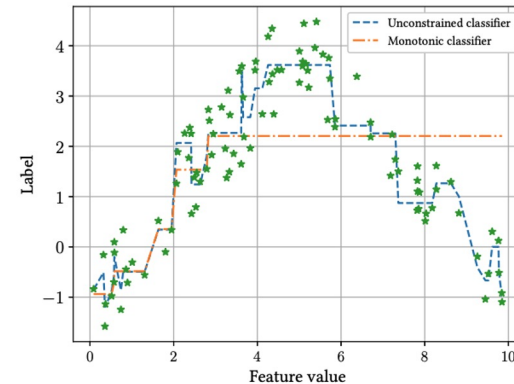
Adversarially Robust Malware Detection Using Monotonic Classification

Inigo Incer*
UC Berkeley
inigo@eecs.berkeley.edu

Sadia Afroz
UC Berkeley
International Computer Science Institute
sadia@icsi.berkeley.edu

Michael Theodorides*
UC Berkeley
theodorides@cs.berkeley.edu

David Wagner
UC Berkeley
daw@cs.berkeley.edu



Certified Detector

Formal guarantees

Proofs of non-existence of adversarial examples around input, leveraging edit distance functions.
Formalized on static end-to-end detectors

Not reproducible

There are no pre-trained available, nor public source code that can be used to train a model

Published

Accepted at NeurIPS 2023!
Also, many work in this direction are appearing in the state of the art

Certified Robustness of Learning-based Static Malware Detectors

Zhuoqun Huang
zhuoqun@unimelb.edu.au
University of Melbourne
Parkville, VIC, Australia

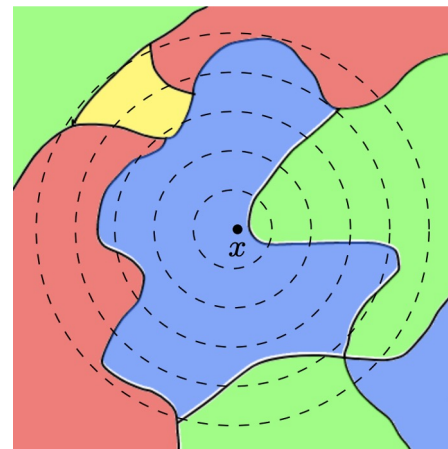
Lujo Bauer
lbauer@cmu.edu
Carnegie Mellon University
Pittsburgh, PA, USA

Neil G. Marchant
nmarchant@unimelb.edu.au
University of Melbourne
Parkville, VIC, Australia

Olga Ohrimenko
oohrimenko@unimelb.edu.au
University of Melbourne
Parkville, VIC, Australia

Keane Lucas
keanelucas@cmu.edu
Carnegie Mellon University
Pittsburgh, PA, USA

Benjamin I. P. Rubinstein
brubinstein@unimelb.edu.au
University of Melbourne
Parkville, VIC, Australia



Upcoming innovation in certification

Only padding?

Recent work only certifies against padding and content-editing attacks (easier to formalize)

Incoming new certification

Joint work on certification for also content-injection attacks that breaks the current methodologies
(combined effort with Daniel Gibert!)

Specific chunking system

Divide incoming input into chunk according to the format, independently from the size of the file or a fixed number of windows.

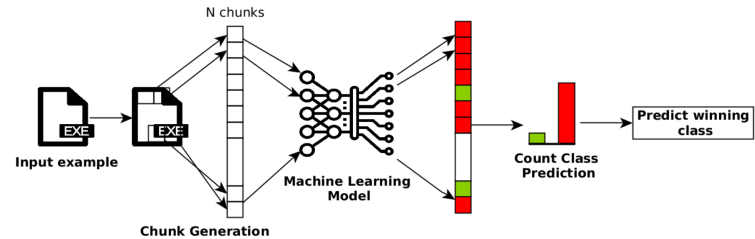
Intuition: the adversarial content will always be contained in contiguous blocks, since it must be aligned with a specific entry in the PE file format

Certified Robustness of Static Deep Learning-based Malware Detectors against Patch and Append Attacks

Daniel Gibert
CeADAR, University College Dublin
Dublin, Ireland
daniel.gibert@ucd.ie

Giulio Zizzo
IBM Research Europe
Dublin, Ireland
giulio.zizzo2@ibm.com

Quan Le
CeADAR, University College Dublin
Dublin, Ireland
quan.le@ucd.ie



Detectors of EXEmples

Aiding Avs without removing them

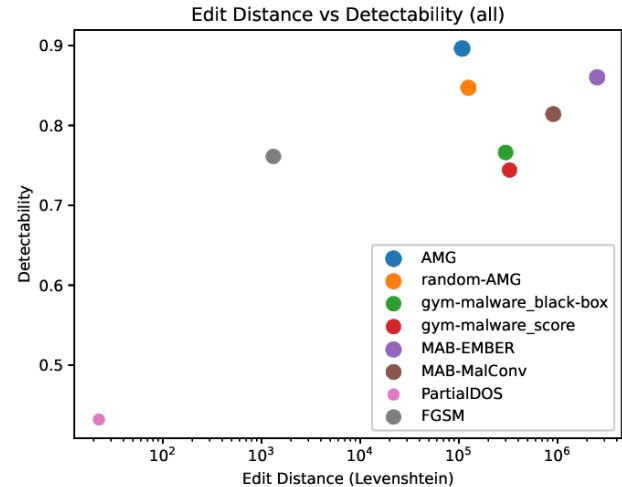
Since it might be difficult to replace a distributed model, we are working to develop a plugin that detects the presence of anomalous EXEmples in test data

Perturbation-spotting

Upcoming work that will show how samples can be discarded if labelled as EXEmples, minimal computational overhead and reduced false positives (thank you Matous Kozak, Ph.D. student @ CTU for this work!)

Trade-off between stealth / effective

The least the EXEmples is modified, the least is detected by scanners: it is a naïve finding, but end-to-end models are more susceptible to “invisible” manipulations than other models!



Take-home messages of Part 3

Research is taking flight

First there was only pre-processing, now we have adversarial training and certification

Need for responsible evaluations

We want to avoid the same history of adversarial defenses on vision models; these new detectors must be evaluated following the practices developed so far

Part 4: Limitations (and future work)

Three big issues with Adversarial EXEmples

Manipulations are hard to craft

Lack of documentation,
lack of open-source reference code,
lack of easily-deployable
debugging tools, and tons of hours
to work

Few available detectors to test

Academic models are either not working
or preliminary, while commercial models
are unavailable

Evasion is not robustness

In system security, evasion should be
achieved “no matter what”,
which is not the same as adversarial
robustness

Creating practical manipulations is painful

Format is vague

Microsoft released vague documentations for the internal of the Windows OS, and research is done by reverse engineering

Debugging is hard

Manipulations often deal with very specific steps of the loader or runtime execution and it is difficult to get messages from the OS

PE Format

Article • 06/23/2022 • 127 minutes to read • 15 contributors



This specification describes the structure of executable (image) files and object files under the Windows family of operating systems. These files are referred to as Portable Executable (PE) and Common Object File Format (COFF) files, respectively.

Note

This document is provided to aid in the development of tools and applications for Windows but is not guaranteed to be a complete specification in all respects. Microsoft reserves the right to alter this document without notice.

This app can't run on your PC

To find a version for your PC, check with the software publisher.

Close

Format specifications are not specific at all

Many details are omitted

Microsoft released an official format documentation, but it is not complete (as they clearly state)

Example: some header fields are not used by the loader, but they are described as meaningful

Windows loader changes through time

It has been proven that Windows XP, 7, and 10 have different loaders that parse the PE structure in a different way!

Closed-source code is not helping

No reference and no code: the only way is either test manipulation by hand, or develop complex tools that infer information about constraints

Note

This document is provided to aid in the development of tools and applications for Windows but is not guaranteed to be a complete specification in all respects. Microsoft reserves the right to alter this document without notice.

Lost in the Loader: The Many Faces of the Windows PE File Format

Dario Nisi
EURECOM
dario.nisi@eurecom.fr

Yanick Fratantonio
Cisco Talos
yfratant@cisco.com

Mariano Graziano
Cisco Talos
magrazia@cisco.com

Davide Balzarotti
EURECOM
davide.balzarotti@eurecom.fr

	Discrepancies				
	W1	W2	W3	W4	W5
XP vs 7	✓	✓			✓
XP vs 10		✓	✓	✓	
7 vs XP					
7 vs 10			✓	✓	
10 vs XP					
10 vs 7	✓				✓

Complex pipeline for debugging manipulations

Dealing with kernel components

The building blocks of the operating systems are inside the kernel, there is no easy way to connect them to a debugger.
Work-around: manual inspection and tons of wasted hours

The only output is the error

Perturbed sample is not working? Keep digging without any other informative log, or rely on other PE viewer or checker (PE Bear, PE Explorer, LIEF, pefile...)

No constraints check

Utilities are good, but they do not tell you IF there is a format specification problem, or they signal vague alerts (if you are lucky)

The screenshot displays the PE Bear v0.3.7 interface. The top window shows the PE structure for 'Locky.exe' and 'locky_25Feb', including sections like '.text', '.data', and '.rsrc'. Below this is a table of sections with columns for Name, Raw Addr, Raw size, Virtual Addr, Virtual Size, Characteristics, Ptr to Reloc, Num. of Reloc., and Num. of Linenum.

The bottom window shows the disassembly of the selected section. The assembly code includes instructions such as:

```
0101F1F9 8B30E20000 push 00000220h
0101F1FE 50 push esp
0101F1FF E805402000 call !_imp_SHELL32.dll!ShellExecuteExW
0101F204 85C9 test ecx,ecx
0101F205 745E jz !L0011F20B
0101F208 800C0001 push byte ptr [!L0011F208+1]
0101F209 8B30E20000 push 00000220h
0101F20C 746F01000000 mov dword ptr [ebp+04h],00000000h
0101F20F FF3540200001 call !_imp_KERNEL32.dll!GetModuleHandleW
0101F214 50 push esp
0101F215 40 call !_imp_KERNEL32.dll!GetProcAddress
0101F218 50 push esp
0101F219 6403 push 00000000h
```

The status bar at the bottom indicates the current cursor position: '39004 EP: 0101F26h Ready... 00:00:00 20:28:50 30.12.2008'.

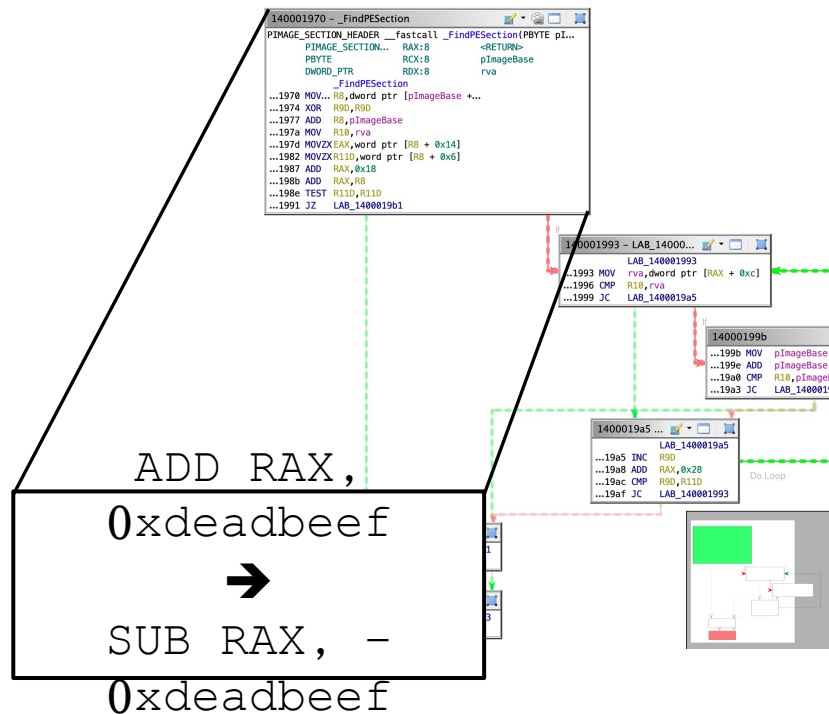
What about attacks against dynamic classifiers?

Not only the structure, the code too!

Code is the only structure that can be manipulated now, so the attacker must act accordingly: code re-writing techniques!

Different instructions, same functionality

Code is re-written to satisfy properties that the attacker wants, like adding new API calls, invert IF statements, add never-to-be executed code to obfuscate...



In practice: a nightmare scenario!

Problems with addresses

If not correctly handled, content injection will shift all known offsets that the compiler created at compile-time

Problems with executable sections

One could create jumps to other code sections... if they are flagged as executable!

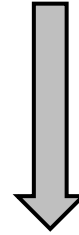
Problems with relocations

Program were usually loaded in the same virtual space, but not secure! The OS randomizes the addresses... and this implies that also adversarial content must take this randomization into account!

Morale: harder than before

It is doable, as there is tons of tools that obfuscate and pack samples, but debugging time levitates from a few to many more hours of human work

```
0x7800: MOV EDI, 1
0x7804: MOV ESI, 2
0x7808 CALL MY_FUNC
. . .
. . .
0x7880 MY_FUNC
```



```
0x7800: MOV EDI, 1
0x7804: MOV ESI, 2
0x7808: XOR EAX, EAX
0x780a CALL MY_FUNC
. . .
. . .
0x7884 MY_FUNC
```

Call for action (1/3)

Develop attacks against DYNAMIC classifiers

Very little has been done, literally two papers that are not reproducible right now
Working on this topic will be key point in testing ALL machine learning malware detectors

Hard? ABSOLUTELY

Rewarding? ABSOLUTELY

Three big issues with Adversarial EXEmples

Manipulations are hard to craft

Lack of documentation,
lack of open-source reference code,
lack of easily-deployable
debugging tools, and tons of hours
to work

Few available detectors to test

Academic models are either not working
or preliminary, while commercial models
are unavailable

Evasion is not robustness

In system security, evasion should be
achieved “no matter what”,
which is not the same as adversarial
robustness

Academic classifiers are (kinda) flawed

Academic models are all “unpublished”

Except for MalConv¹, all the other classifiers used in research are published as preprint: EMBER, PEberus, Non-Negative MalConv², and many others

No adversarial robustness is considered

Except for PEberus, no model consider adversarial attacks inside their formulation, and they can be easily evaded

Only static detectors

There are no open-source state-of-the-art machine learning models that rely on runtime information of Windows programs³

1. MalConv can be evaded by replacing ~60 bytes
2. Non-Negative MalConv has an open-source release that has a terrible ROC

Against All Odds: Winning the Defense Challenge in an Evasion Competition with Diversification

This repository contains the defense *PEberus** that got the first place in the [Machine Learning Security Evasion Competition 2020](#), resisting a variety of attacks from independent attackers.

You can find the whitepaper that outlines our defense here:

Elastic Malware Benchmark for Empowering Researchers

The EMBER dataset is a collection of features from PE files that serve as a benchmark dataset for researchers. The EMBER2017 dataset contained features from 1.1 million PE files scanned in or before 2017 and the EMBER2018 dataset contains features from 1 million PE files scanned in or before 2018. This repository makes it easy to reproducibly train the benchmark models, extend the provided feature set, or classify new PE files with the benchmark models.

This paper describes many more details about the dataset: <https://arxiv.org/abs/1804.04637>



Industry classifiers are locked away

Commercial means Unavailable

All the companies do not share any technical insights about their technologies (of course), most of them can not be tested with a free license, other must be reverse engineered

VirusTotal is the only way

Most of them can be tested using the crowd service VirusTotal, but only old not-updated versions are available



These vendors sell machine learning inside their products.

Call for action (2/3)

Develop new state-of-the-art models

We are still using models from 2018, more focus on modelling defenses than creating good and easy-to-use models like MalConv and GBDT

Three big issues with Adversarial EXEmples

Manipulations are hard to craft

Lack of documentation,
lack of open-source reference code,
lack of easily-deployable
debugging tools, and tons of hours
to work

Few available detectors to test

Academic models are either not working
or preliminary, while commercial models
are unavailable

Evasion is not robustness

In system security, evasion should be
achieved “no matter what”,
which is not the same as adversarial
robustness

Evasion is not equal to adversarial robustness

Two different goals

Evasion implies that malware samples bypass detection, while adversarial robustness quantifies the sensitivity of the detector

Interested in evasion? Obfuscate & Pack

Tons of literature, open-source code, and material to evade ANY malware detector (with or without machine learning). Since tools are automatic, it does not change much to perturb or inject few bytes or kilobytes

Static detection is bypassed by design

Structure of programs can be changed and embedded in other programs, downloaded from the internet after execution, and more



MPRESS

Version: 2.19



Very well known packer programs that hide malicious content. Originally created to prevent reverse engineering of legitimate code.

Do companies care about Adversarial robustness?

“We really appreciate this research and would like to collaborate to continue to improve our products and services. At this time this technique does not meet the definition of vulnerability in the product or has demonstrated that it bypassed our products. We will however look into our static ML models to see how we can incorporate this technique to further improve.”

One-of-those-company

Adversarial EXamples not treated as vulnerability

Companies are interested in evading the overall pipeline, not just portions of the products

Naïve solution: test attacks against deployed commercial products

(which are unavailable, as said before)

Joshua Saxe
@joshua_saxe

Why robustness to adversarial examples isn't a first-priority concern on the Sophos AI team.

Why adversarial examples aren't a primary concern at Sophos AI

The decision function the adversarial ML literature assumes we use for PE file detection

A real-world decision function where user downloads and executes a PE file from the web

Check if PE ML model detects file

Frequently updated URL, blocklist, heuristic, signature, and ML checks

Frequently updated HTML/CSS blocklist/allowlist, signature and ML checks

Downloaded PE file signature, blocklist, allowlist and ML checks

Blocking post-execution behavioral detection checks (rule-based)

Non-blocking post-execution behavioral rule and ML alerting, for follow-up by SOC teams

The adversarial example research frame makes perfect sense here

Here the adversary is attempting to bypass many non-ML layers that change dynamically over time, and that block it on reconnaissance of the whole system.

8:17 PM · Jul 22, 2022 · Twitter Web App

Missing the bigger picture

Companies have the feelings that academic settings are unrealistic, as they target “only” the ML component

Call for action (3/3)

Develop new testing techniques that consider all components

We are starting to create end-to-end pipelines, but we are still missing a complete framework that systematically tell a developer HOW to test these models

Take-home messages of Part 4

Manipulation are hard to craft

Requires patience and hours of work, high risk / high reward scenario

Few classifiers around

We are still using EMBER from 2018 (5 years ago!) with no candidate that performs better, both in terms of accuracy and robustness

Industry are not so concerned (yet)

The focus is mostly shifted towards controlling false positives

Thanks!



Luca Demetrio

luca.demetrio@unige.it

X @zangobot



*If you know the enemy and know yourself, you need not fear
the result of a hundred battles*

Sun Tzu, The art of war, 500 BC