

# Messing with binary formats

ΠΟΛΥΓΛΩΣΣΟΣ

44CON



Ange Albertini  
2013/09/13

London, England



# Welcome!

- this is the non-live version of my slides
  - more text
  - standard PDF file ;)

## About me:

- Reverse engineer
- my website: <http://corkami.com>
  - reverse engineering & visual documentations



**I just like to play  
with lego blocks**



low-level ones,  
that is





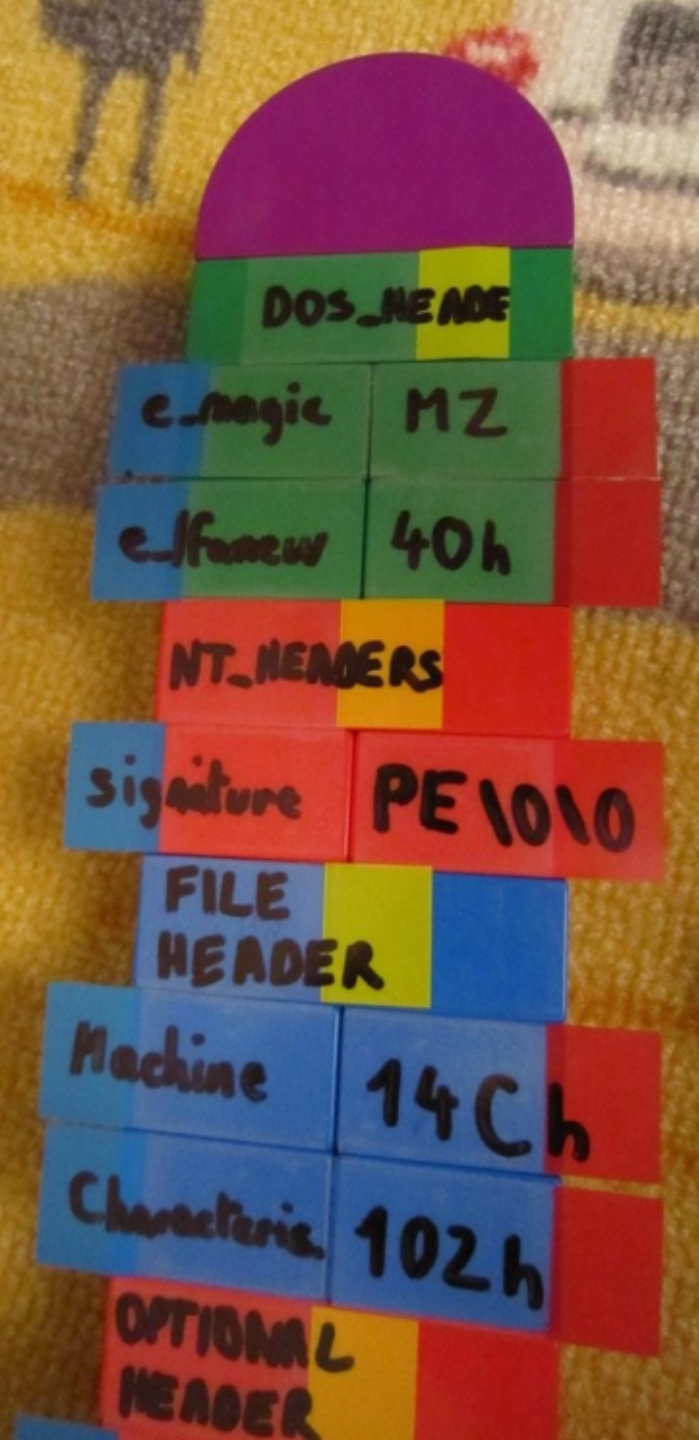
# block by block

## generate files byte per byte

```
istruc IMAGE_DOS_HEADER
    at IMAGE_DOS_HEADER.e_magic, db 'MZ'
    at IMAGE_DOS_HEADER.e_lfanew, dd NT_Signature - IMAGEBASE
iend
NT_Signature:
istruc IMAGE_NT_HEADERS
    at IMAGE_NT_HEADERS.Signature, db 'PE', 0, 0
iend
istruc IMAGE_FILE_HEADER
    at IMAGE_FILE_HEADER.Machine, dw IMAGE_FILE_MACHINE_I386
    at IMAGE_FILE_HEADER.Characteristics, dw IMAGE_FILE_EXECUTABLE_IMAGE
iend
istruc IMAGE_OPTIONAL_HEADER32
```

## Goals

- explore the format
- make sure that's how things work
- full control over the structure



# a complete executable

IMAGEBASE equ 400000h

org IMAGEBASE

istruc IMAGE\_DOS\_HEADER

at IMAGE\_DOS\_HEADER.e\_magic, db 'MZ'

at IMAGE\_DOS\_HEADER.e\_lfanew, dd NT\_Signature - IMAGEBASE

iend

NT\_Signature:

istruc IMAGE\_NT\_HEADERS

at IMAGE\_NT\_HEADERS.Signature, db 'PE', 0, 0

iend

istruc IMAGE\_FILE\_HEADER

at IMAGE\_FILE\_HEADER.Machine, dw IMAGE\_FILE\_MACHINE\_I386

at IMAGE\_FILE\_HEADER.Characteristics, dw IMAGE\_FILE\_EXECUTABLE\_IMAGE

iend

istruc IMAGE\_OPTIONAL\_HEADER32

at IMAGE\_OPTIONAL\_HEADER32.Magic, dw IMAGE\_NT\_OPTIONAL\_HDR32\_MAGIC

at IMAGE\_OPTIONAL\_HEADER32.AddressOfEntryPoint, dd EntryPoint - IMAGEBASE ; not strictly required

at IMAGE\_OPTIONAL\_HEADER32.ImageBase, dd IMAGEBASE ; not required under XP

at IMAGE\_OPTIONAL\_HEADER32.SectionAlignment, dd 1

at IMAGE\_OPTIONAL\_HEADER32.FileAlignment, dd 1

at IMAGE\_OPTIONAL\_HEADER32.MajorSubsystemVersion, dw 4

at IMAGE\_OPTIONAL\_HEADER32.SizeOfImage, dd SIZEOFIMAGE

at IMAGE\_OPTIONAL\_HEADER32.SizeOfHeaders, dd SIZEOFIMAGE - 1 ; required for XP

at IMAGE\_OPTIONAL\_HEADER32.Subsystem, dw IMAGE\_SUBSYSTEM\_WINDOWS\_CUI

iend

istruc IMAGE\_DATA\_DIRECTORY\_16

iend

EntryPoint:

push 42

pop eax

ret

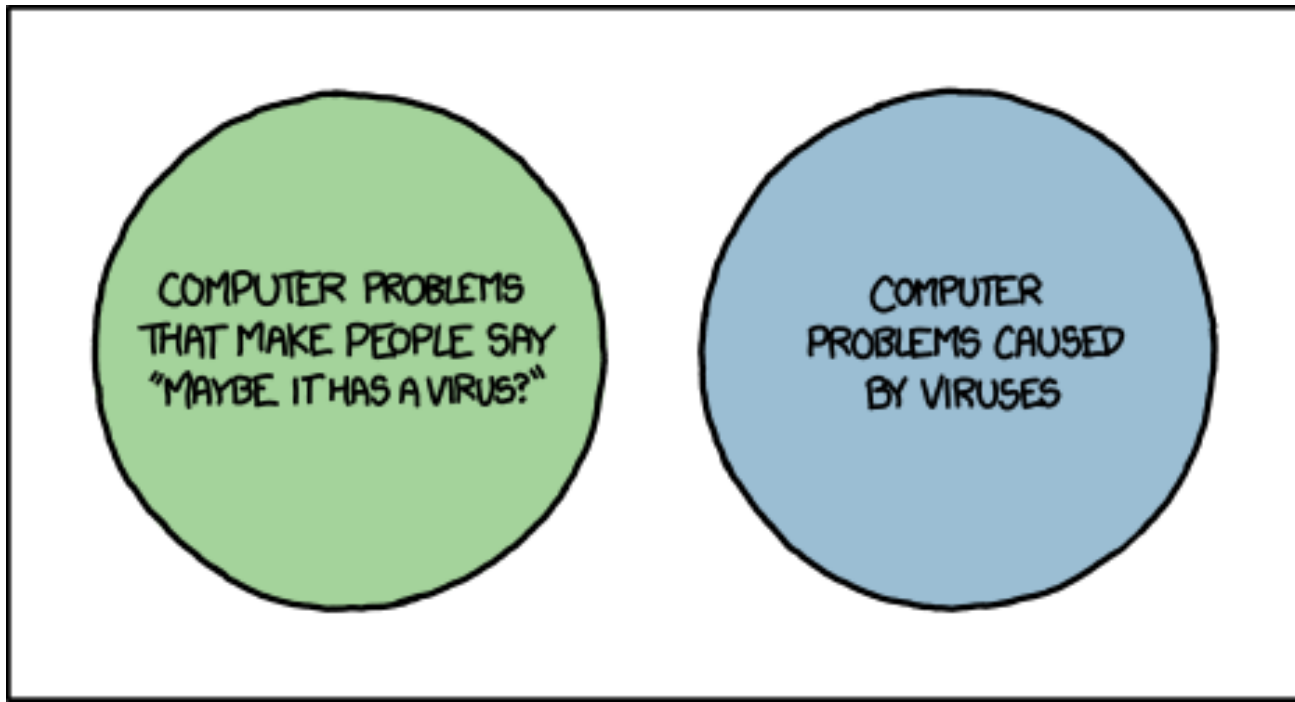
result:

- a complete executable
- all bytes defined by hand



# our problem

- is related to virus (malwares)
- they use many file formats
- it's critical to identify them reliably
  - and to tell whether corrupted or well-formed

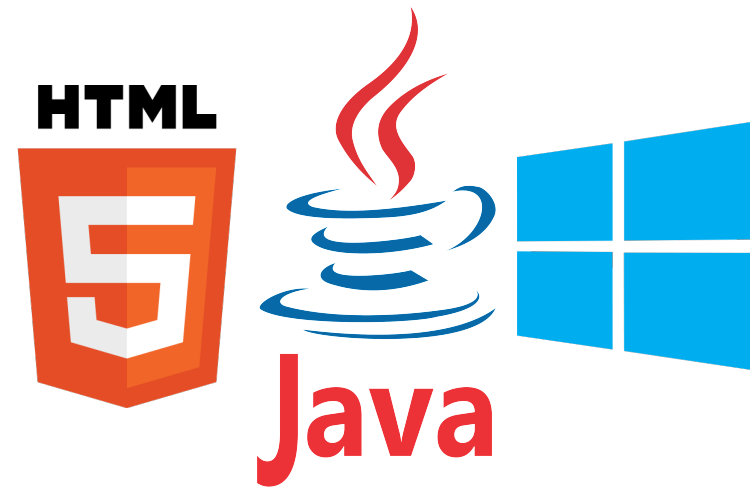


# standard infection chain

the most common chain:

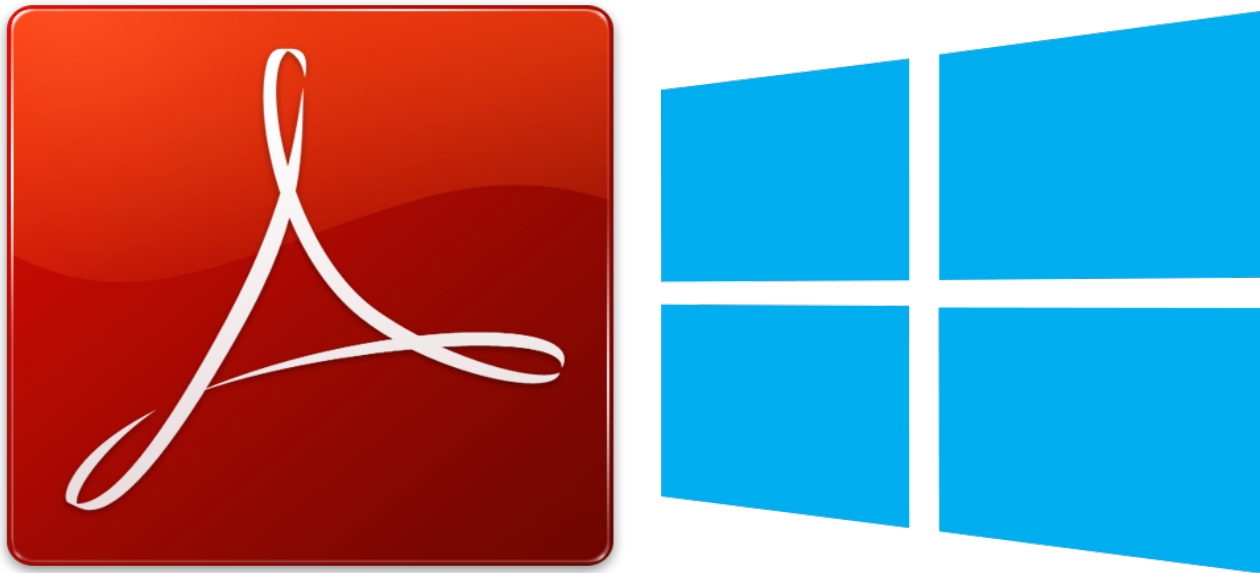
1. a web page, in HTML format
  - a. launching an applet
2. an evil applet, in CLASS format
  - a. exploiting a Java vulnerability
  - b. dropping an executable
3. a malicious executable, in Portable Executable format

(a vast majority of malwares rely on an executable)



# another classic chain

- open a PDF document
  - with an exploit inside
    - dropping or downloading a PE executable
- get a malicious executable on your machine





# the challenge

it might look obvious:

- tell whether it's a PDF, a PE, a JAVA, an HTML...
- typical formats are clearly defined
  - Magic signature enforced at offset 0



# reality

some formats have no header at all

- Command File (DOS 16 bits)
- Master Boot record

some formats don't need to start at offset 0

- Archives (Zips, Rars...)
- HTML
  - but text-only?

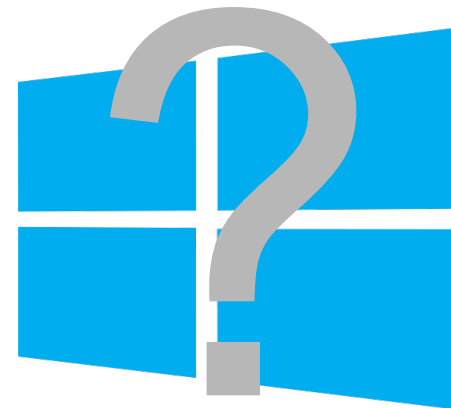
some formats accept a large controllable block early in their header

- Portable Executable
- PICT image

# How did this start?

a real-life problem:

1. a (malicious) HTML page
2. started with 'MZ' (the signature of PE)
3. just scanned as a PE!
  - a. wow, this PE is highly corrupted :)
  - b. it must be clean :p



# polyglots in the wild

GIFAR = GIF + JAR

- an uploaded image
  - an avatar in a forum
- with a malicious JAVA appended as JAR

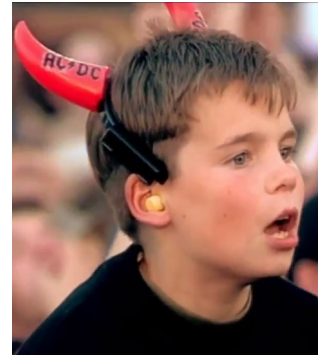
hosted on the server!

- bypass *same domain policy*
- now useable via its *JAVA=EVIL* payload



+ =

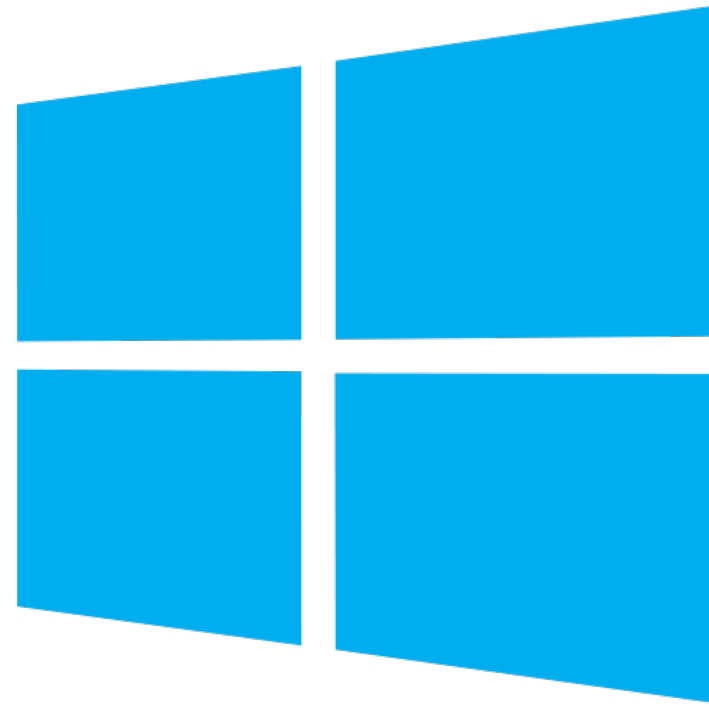
```
<applet  
  code=gifar.class  
  archive=gifar.gif  
>  
</applet>
```



# let's get started

PE, the executable format of windows

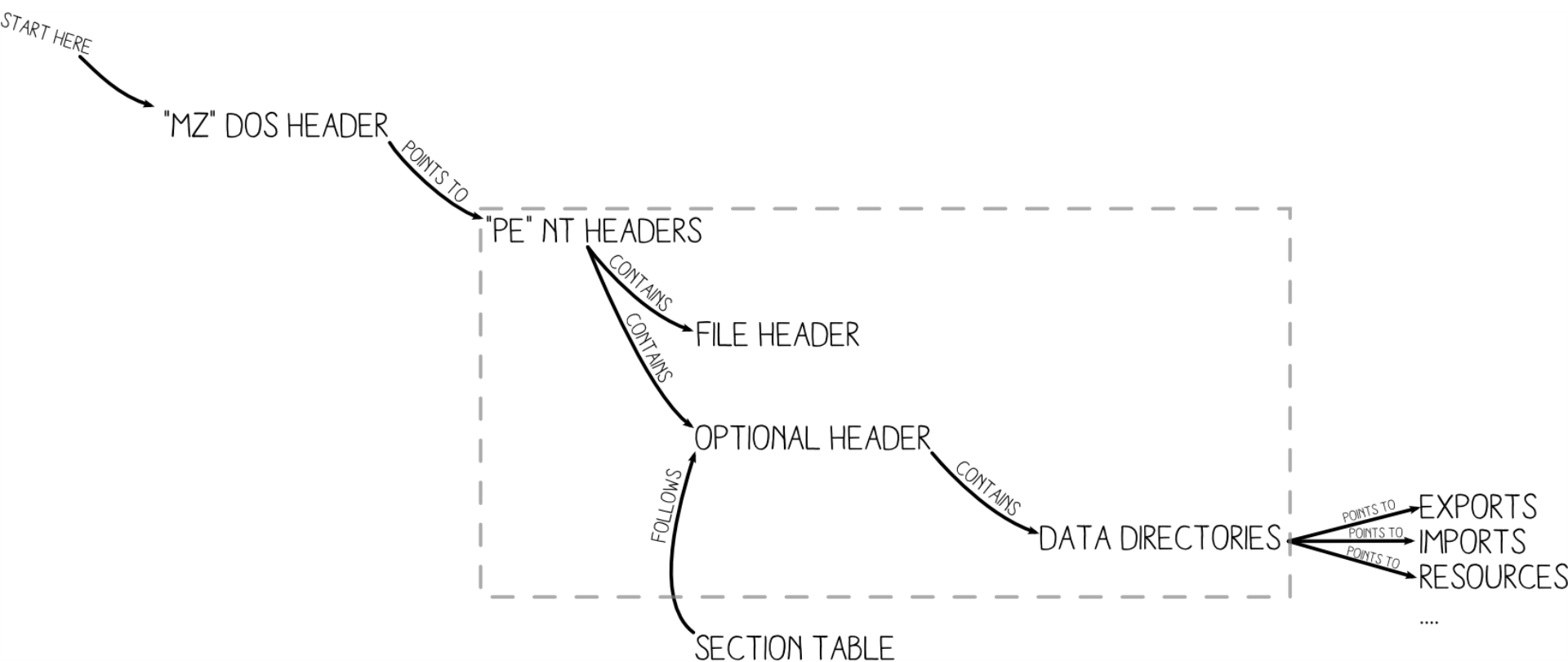
- it's central to windows malware
- it enforces a magic signature at offset 0
  - game over for other formats?





# overview

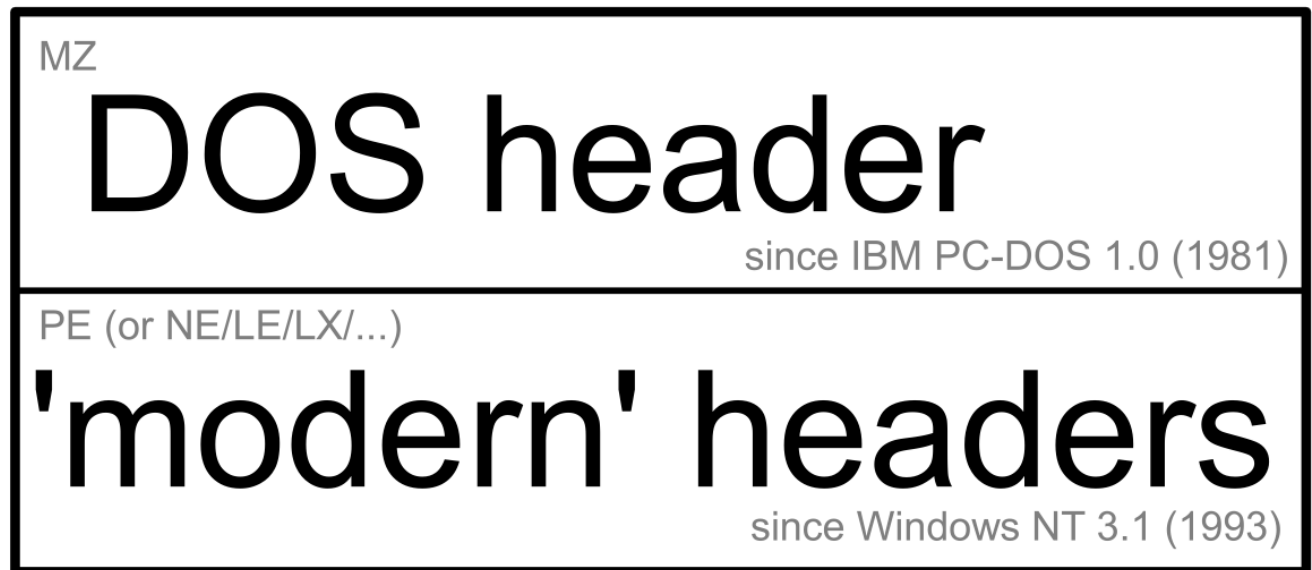
- starts with a compulsory header
- made of sub-headers



# a historical sandwich

1. a deprecated but required header
2. a modern header

## Header



# old header content

- almost completely ignored
- only required:
  - 2 byte signature
  - pointer to new header

OFFSET 0

DOS Header

## IMAGE\_DOS\_HEADER

```
00+2 e_magic MZ
02+2 e_cblp
04+2 e_cp      exe size
06+2 e_crlc
08+2 e_cparhdr exe start
0a+2 e_minalloc
0c+2 e_maxalloc
0e+2 e_ss initial ss
10+2 e_sp initial sp
12+2 e_csum
14+2 e_ip
16+2 e_cs
18+2 e_lfarlc
1a+2 e_ovno
1c+2 e_res[4]
24+2 e_oemid
26+2 e_oeminfo
28+2 e_res2[10]
3c+4 e_lfanew
```

# the new header can be anywhere

ex: at the end of the file!  
such as **Corkami Standard Test**

```
00004A70: 63 6F 64 63 0D 00 43 52 52 4F 52 3A 20 37 36 62 code: ERROR: 760
00004A80: 69 74 20 63-6F 64 65 0D-0A 00 49 6E-66 6F 3A 20 it code: Info:
00004A90: 36 34 20 62-69 74 73 20-6E 6F 74 20-73 75 70 70 64 bits not supp
00004AA0: 6F 72 74 65-64 0D 0A 00-00 00 00 00-00 00 00 00 orted:
00004AB0: 50 45 00 00-4C 01 00 00-D3 F6 5E 81-B1 0F CB 06 PE L0 u: ^ü
00004AC0: 36 06 E7 32-08 01 0F 01-0B 01 8E AF-96 D3 5E A6 6 r2 0 0 0 0 0 0 0 0
00004AD0: ED 48 81 8B-EE CB 6E 38-00 00 00 00-A8 1D DA 96 H i 8 0 0 0 0 0 0
00004AE0: 9B D5 36 CF-00 00 FD 7E-01 00 00 00-01 00 00 00 F6 ± 2 ~ 0
00004AF0: C4 5E A2 35-58 44 C8 EF-04 00 E5 A5-00 00 00 00 - ^ 65 X D L 0 0 0 0
00004B00: D6 4B 00 00-D5 4B 00 00-18 E7 A9 01-03 00 70 9A K fK ↑ r- 0 pü
00004B10: C7 BD 12 00-A8 1A 00 00-06 9E 12 00-23 01 00 00 μ ↑ 0 0 0 0 0 0 0 0
00004B20: A5 2B 4A CE-6D 43 B9 B2-10 26 00 00-68 A8 1A 57 N + J m C i 0 0 0 0 0 0 0 0
00004B30: B8 24 00 00-26 81 CD 27-00 00 00 00-78 EA 0B F5 7 $ 8 ü = ' x Ω δ J
00004B40: 63 0F 2B 56-0C 31 BE 17-8B 67 C2 18-0B 64 F5 D8 c * + V 9 1 d i g T i δ d J †
00004B50: C0 27 00 00-14 00 00 00-08 CA 8D 9A-00 00 00 00 L ' 0 i ü
00004B60: 61 F5 9F CE-CE B3 CF BA-83 3E 46 89-5D 1D 1E F9 a J f i m i | ± | â > F e l + Δ ·
00004B70: FC 00 00 00-02 A1 E7 60-00 00 00 00-E9 88 52 8D " 0 i r 0 ê R i
00004B80: 00 00 00 00-34 DD 53 D4-88 24 00 00-B9 01 00 00 4 | S ê $ | | 0
00004B90: 7C 4F 57 9E-CB D2 98 DC-00 00 00 00-A6 59 CD 93 ! O W R m y 2 Y = ô
00004BA0: E2 D1 5E B4-95 25 BB 0B-FF 35 FC 02-FD 7E E8 2D Γ τ ^ 0 % n δ 5 " 0 2 ~ 0 -
00004BB0: B7 FF FF C3-00 00 00 00-00 00 00 00-00 00 00 00 n |
00004BC0: 0D 0A 00 FF-25 90 24 FD-7E 00 00 00-00 00 00 00 J 0 % é $ 2 ~
00004BD0: FF 25 94 24-FD 7E - - - - - % ö $ 2 ~
1 2GetBlk 3Replac 4ReRead 5Base 6 70ther 8 9FilArg10$avSt
```

**let's look at HTML format**

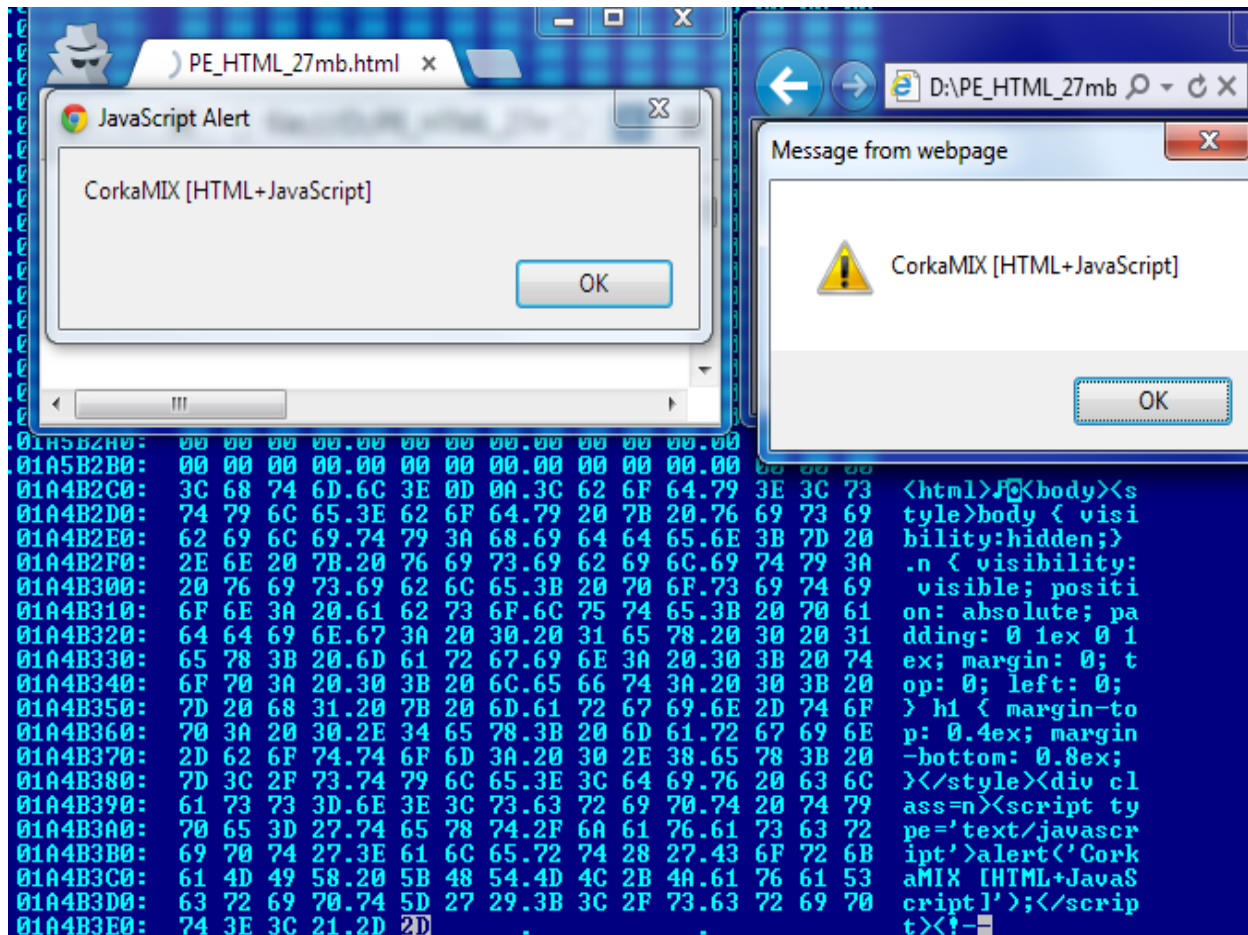
**HTML**





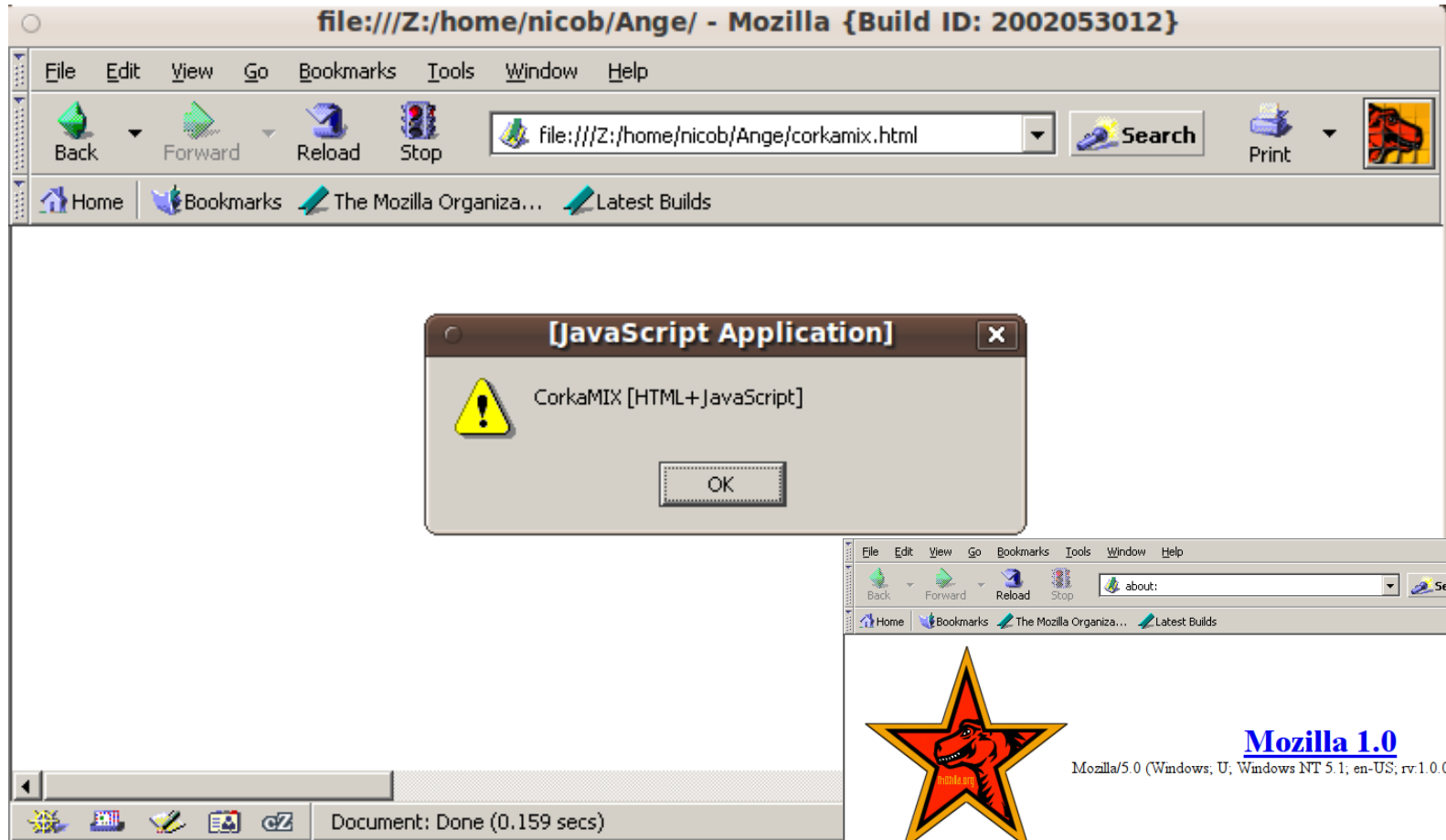
# it enforces NOTHING!

anything before the <html> tag!  
even 28 Mb of binary!



# and it's been the same since **Mozilla 1.0** in 2002

thanks to Nicolas Grégoire!



**now, the PDF format**

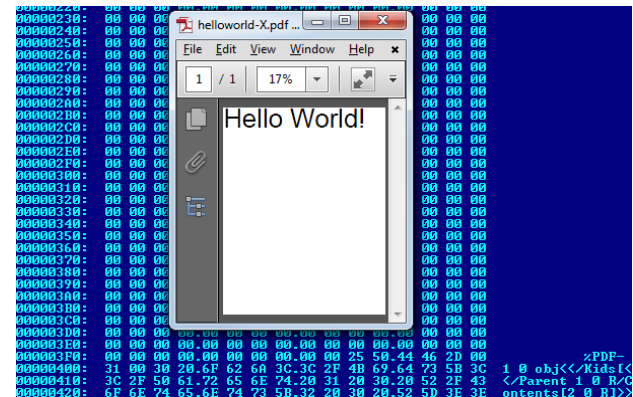


# signature position?

- officially at offset 0
- officially tolerated until offset 1024
- wtf?
  - it get actually worse later

## 7.5.2 File Header

The first line of a PDF file shall be a *header* consisting of the 5 characters `%PDF-` followed by a version number of the form `1.N`, where `N` is a digit between 0 and 7.

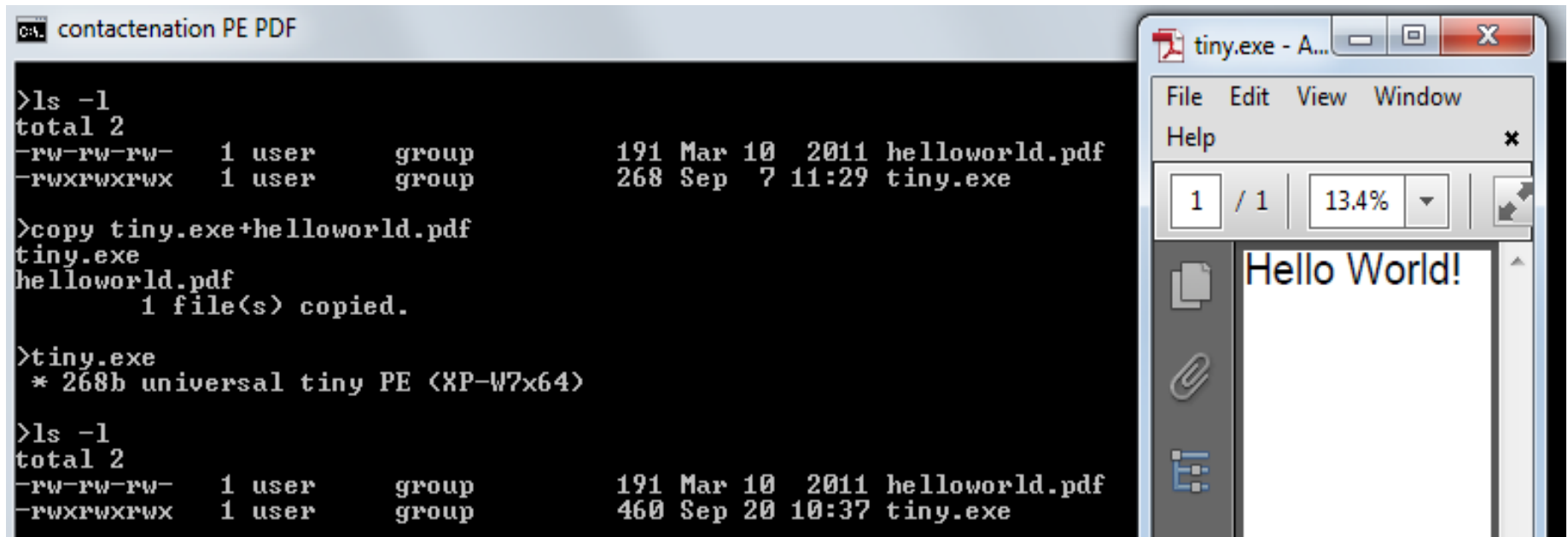


## 3.4.1, "File Header"

13. Acrobat viewers require only that the header appear somewhere within the first 1024 bytes of the file.

# PDF trick 1

put a small executable within 1024 bytes  
(just concatenate)



The image shows a Windows command prompt window titled "contactenation PE PDF" and a small application window titled "tiny.exe - A...".

The command prompt shows the following commands and output:

```
>ls -l
total 2
-rw-rw-rw- 1 user group 191 Mar 10 2011 helloworld.pdf
-rwxrwxrwx 1 user group 268 Sep 7 11:29 tiny.exe

>copy tiny.exe+helloworld.pdf
tiny.exe
helloworld.pdf
1 file(s) copied.

>tiny.exe
* 268b universal tiny PE (XP-W7x64)

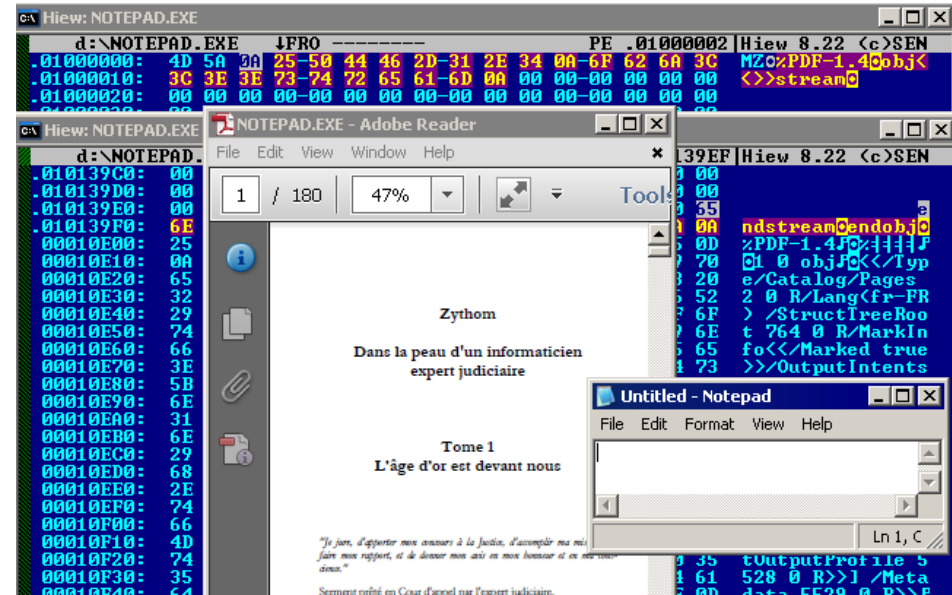
>ls -l
total 2
-rw-rw-rw- 1 user group 191 Mar 10 2011 helloworld.pdf
-rwxrwxrwx 1 user group 460 Sep 20 10:37 tiny.exe
```

The "tiny.exe" window displays the text "Hello World!" in a simple font. The window has a menu bar with "File", "Edit", "View", "Window", and "Help". The status bar shows "1 / 1" and "13.4%".



# trick 2

1. start a fake PDF + object in a PE header
  2. finish fake object at the end the PE
  3. end fake object
  4. put PDF real structure
- works with real-life example!  
(PE data might contain PDF keywords)



# JAR = ZIP + Class

just enforced at the very end of the file



# but CRCs are just ignored

it was too easy :p

```
>yasm -o test.jar zip.asm
```

```
>unzip -lv test.jar
```

Length	Method	Size	Ratio	Date	Time	CRC-32	Name
0	Stored	0	0%	00/00/80	00:00	00000000	META-INF/
35	Stored	35	0%	00/00/80	00:00	deadbeef	META-INF/MANIFEST.MF
299	Stored	299	0%	00/00/80	00:00	0badbabe	test.class
334		334	0%				3 files

```
>unzip -t test.jar
```

```
Archive: test.jar
  testing: META-INF/ OK
  testing: META-INF/MANIFEST.MF bad CRC 8391c53a (should be deadbeef)
  testing: test.class bad CRC 7846a510 (should be 0badbabe)
At least one error was detected in test.jar.
```

```
>java -jar test.jar
Java: Working! (with wrong CRCs)
```

```
>
```

# Summary

# Structure

## 1. start

- PE Signature
  - %PDF + fake *obj* start
  - HTML comment start

## 2. next

- PE (next)
- HTML
- PDF (next)

## 3. bottom

- ZIP



it's time for a real example!

an *inception* demo!

wait, what?



# we're already in the demo!

the **live** version file is simultaneously:

- the PDF slides themselves
- a PDF viewer executable
  - ie, the file is loading itself
- the PoCs in a ZIP
- an HTML readme
  - with JavaScript mario



# **so, it works**

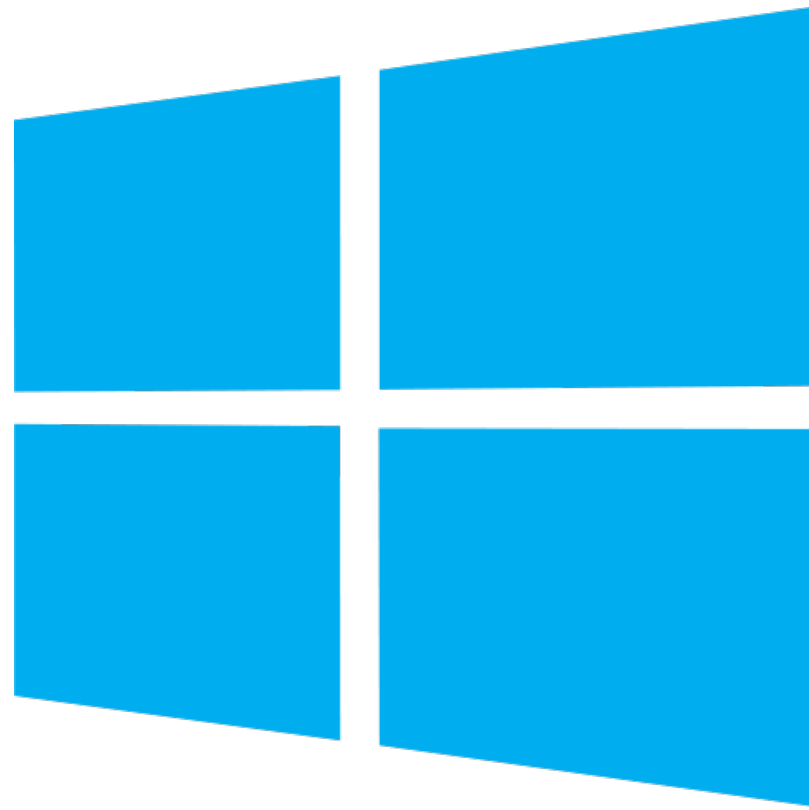
but it lacks something

- not artistic enough
- not advanced enough

let's build a 'well representative' (=nasty) PoC

# the PE specs

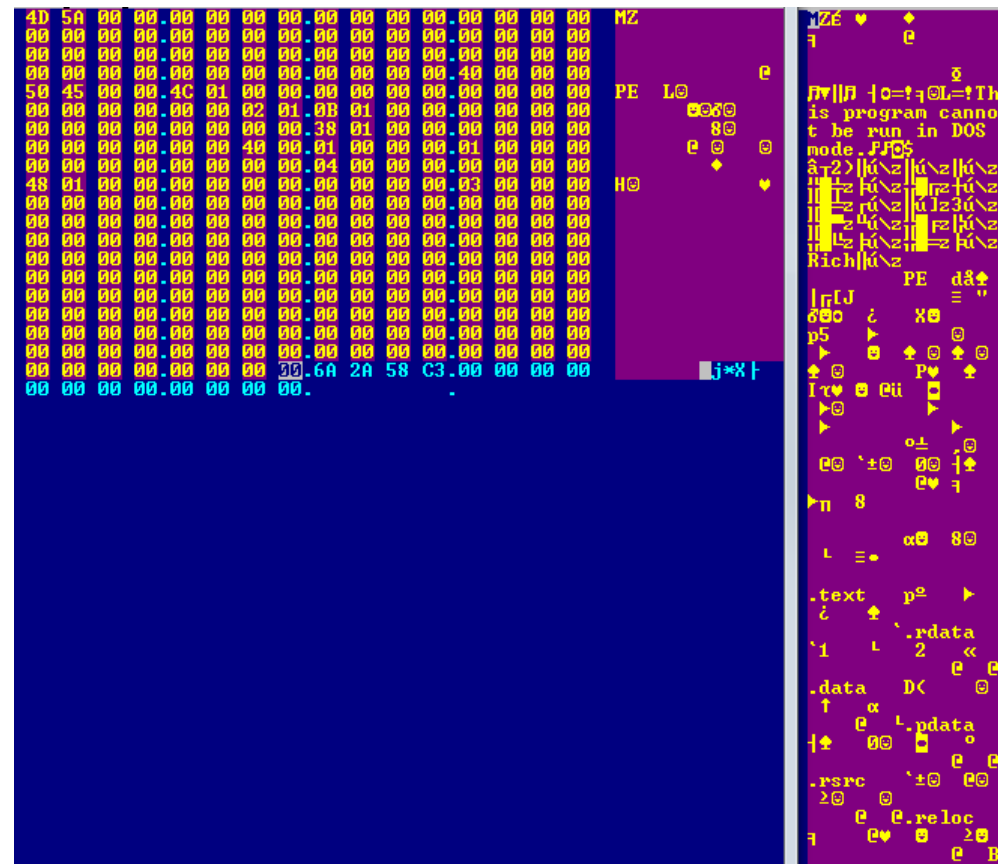
- Official MS specs = big joke
  - 'the gentle guide for beginners'
  - barely describes standard PEs



# stripped down PE

many elements removed

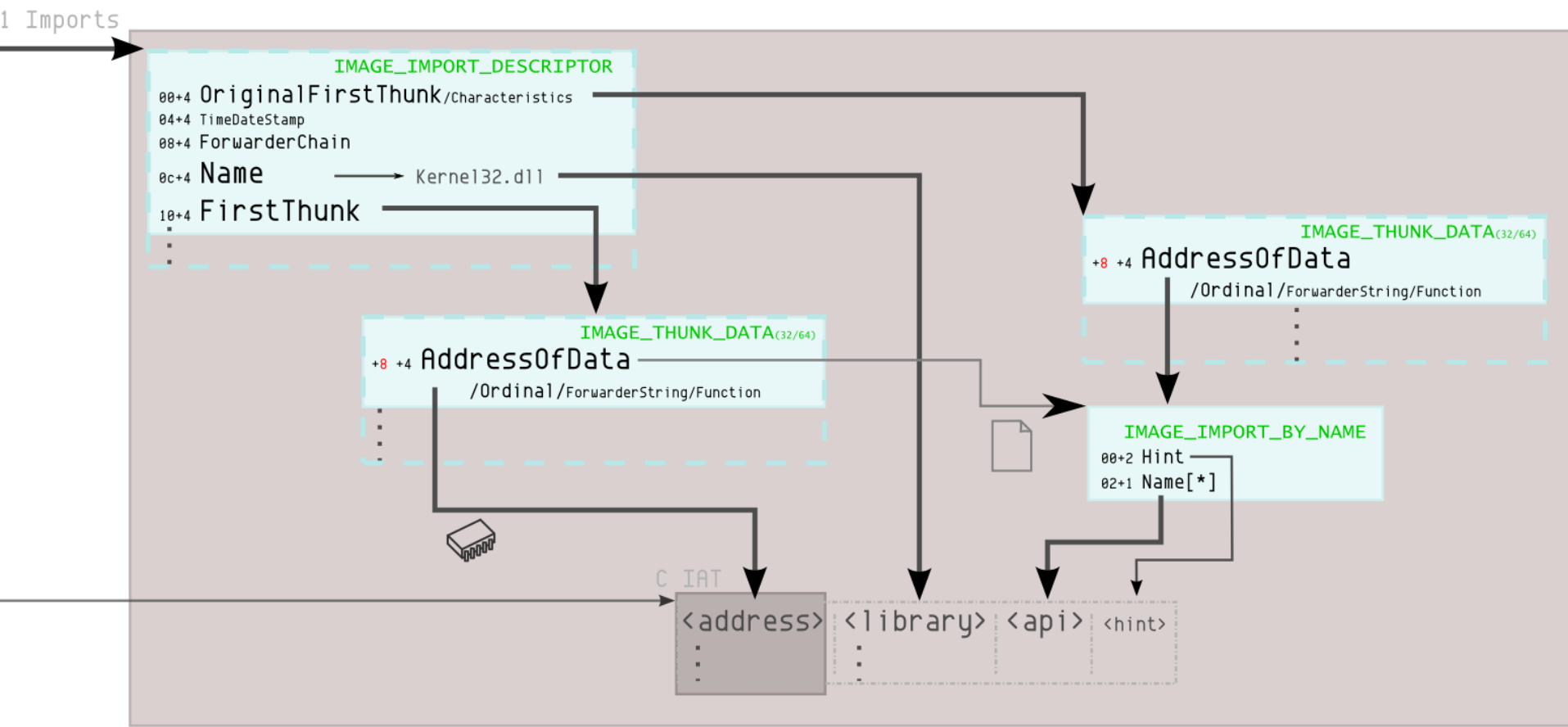
- including no sections



# imports

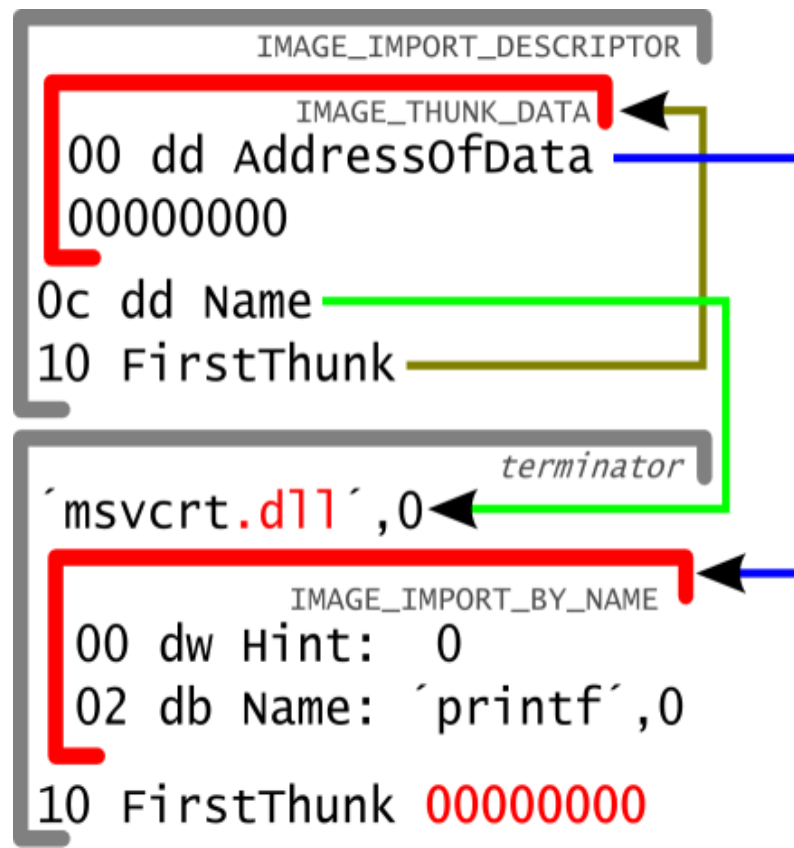
(imports = communication between executables and libraries)

imports are made of 3 lists



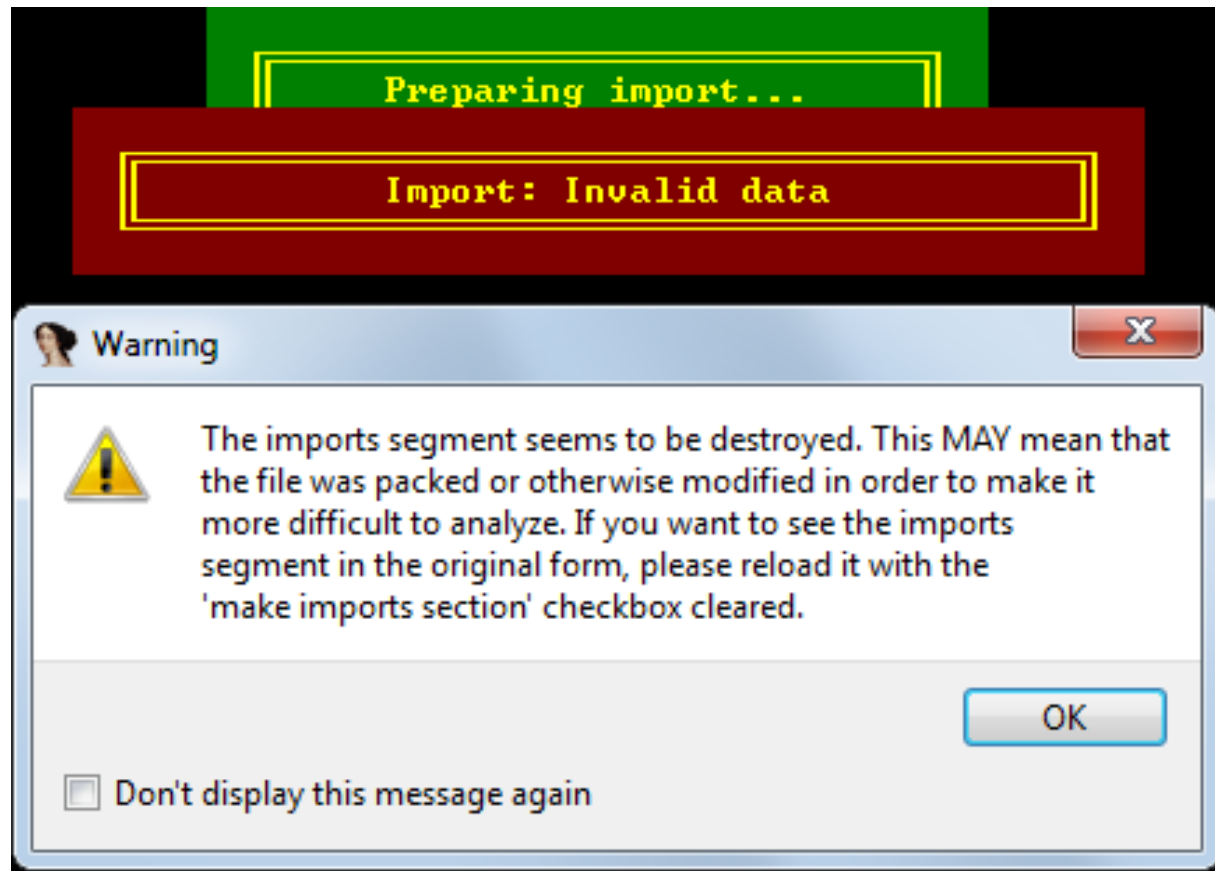
# evil imports

- let's make these lists into each other
- with more extra tricks to fail parser!



# ultimate import fail

- failing all tools
  - including IDA & Hiew
- now fixed :)





# let's put some code

- some undocumented opcodes!
- big blank spaces in Intel official docs

Table A-2. One-byte Opcode Map: (00H — F7H) \*

	0	1	2	3	4	5	6	7
D	Eb, 1	Ev, 1	Eb, CL	Ev, CL	AAM <sup>i64</sup> lb	AAD <sup>i64</sup> lb		XLAT/ XLATB

Table A-3. Two-byte Opcode Map: 08H — 7FH (First Byte is 0FH) \*

	px	8	9	A	B	C	D	E	F
0		INVD	WBINVD		2-byte Illegal Opcodes UD2 <sup>1B</sup>		NOP Ev		
1		Prefetch <sup>1C</sup> (Grp 16 <sup>1A</sup> )							NOP Ev
		vmmovs	vmmovs	cvtpi2ps	vmmovps	cvtps2pi	cvtps2pi	vmmovss	vmmovss

## let's check AMD's

- miracle!

### Table A-1. One-Byte Opcodes, Low Nibble 0–7h

Nibble <sup>1</sup>	0	1	2	3	4	5	6	7
D	Eb, 1	Ev, 1	Eb, CL	Ev, CL	AAM <sup>3</sup>	AAD <sup>3</sup>	SALC <sup>3</sup>	XLAT

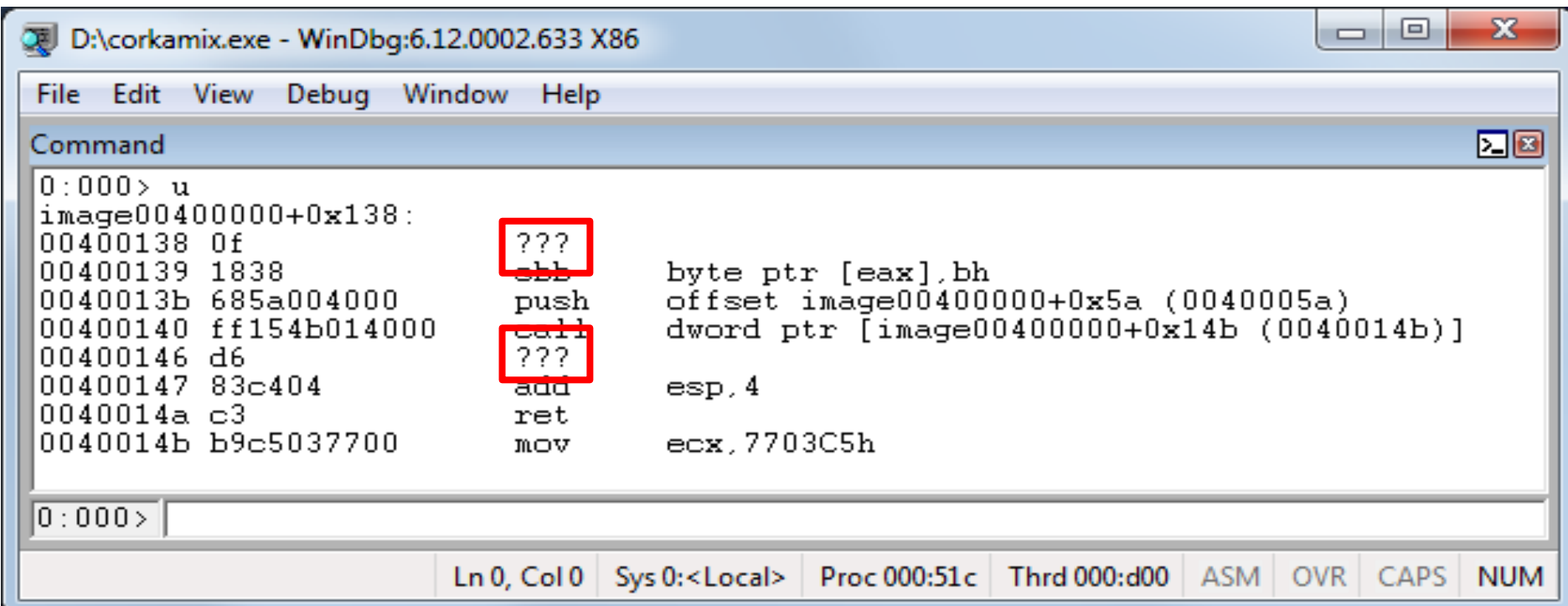
### Table A-4. Second Byte of Two-Byte Opcodes, Low Nibble 8–Fh

[illegible]

# result in WinDbg

- '???' == clueless (tool/user)

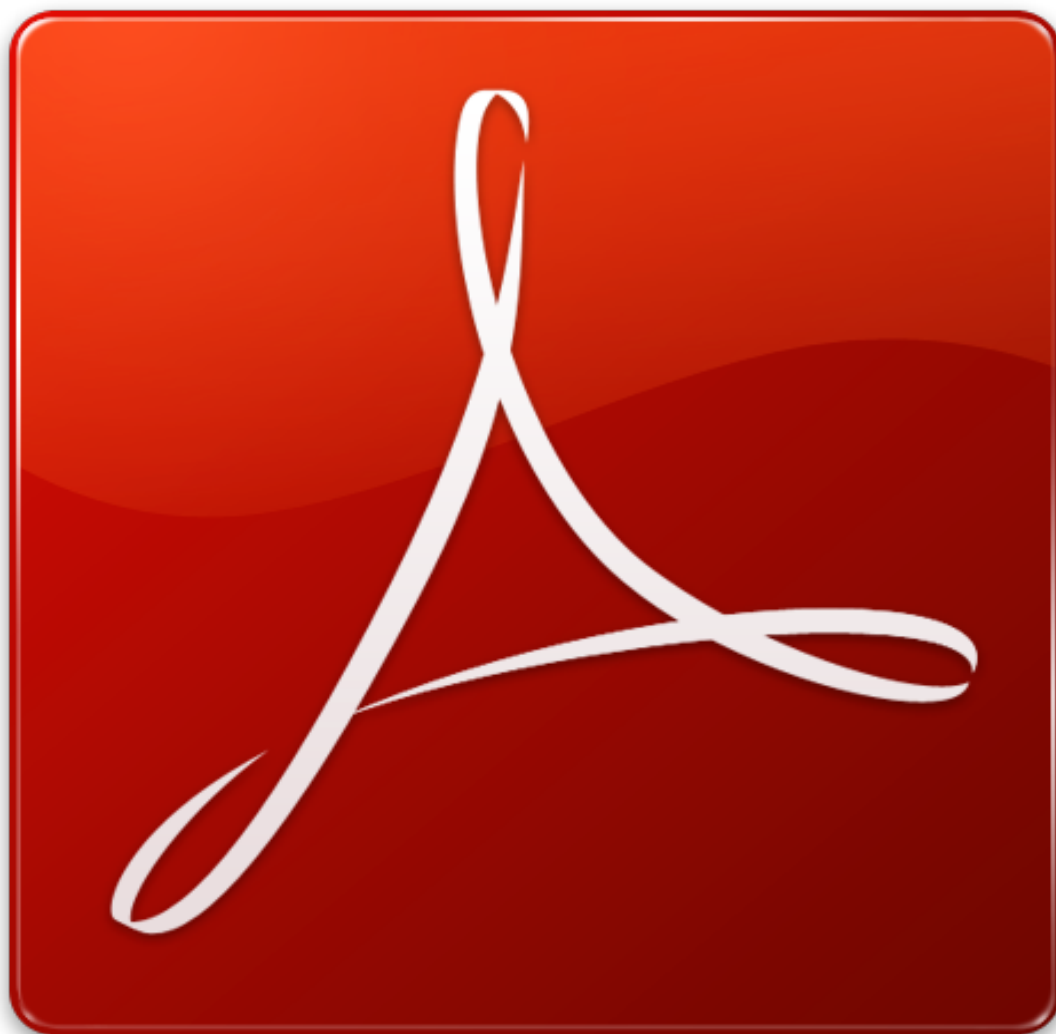
**don't rely (only) on official docs**



```
D:\corkamix.exe - WinDbg:6.12.0002.633 X86
File Edit View Debug Window Help
Command
0:000> u
image00400000+0x138:
00400138 0f          ???
00400139 1838       cbb      byte ptr [eax],bh
0040013b 685a004000 push     offset image00400000+0x5a (0040005a)
00400140 ff154b014000 call     dword ptr [image00400000+0x14b (0040014b)]
00400146 d6        ???
00400147 83c404     add     esp,4
0040014a c3        ret
0040014b b9c5037700 mov     ecx,7703C5h
0:000>
```

Ln 0, Col 0   Sys 0:<Local>   Proc 000:51c   Thrd 000:d00   ASM   OVR   CAPS   NUM

# **messing with PDF**

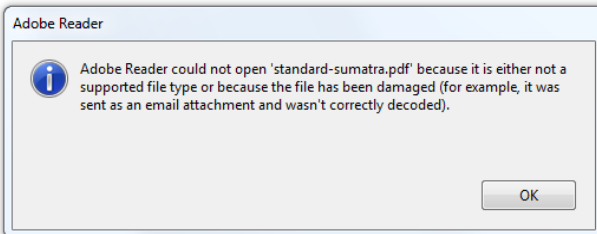
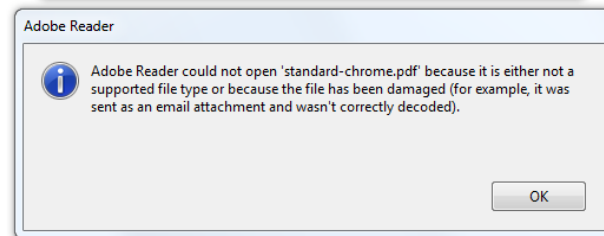
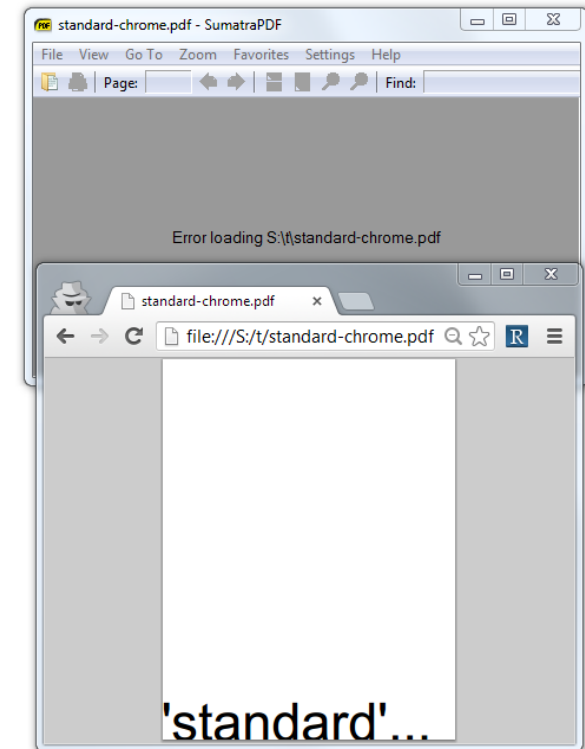
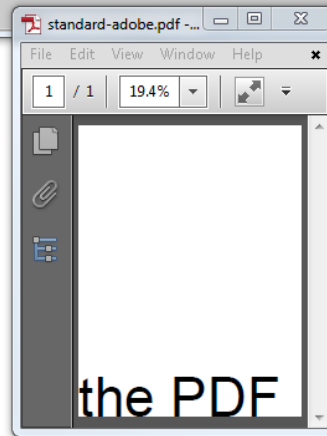
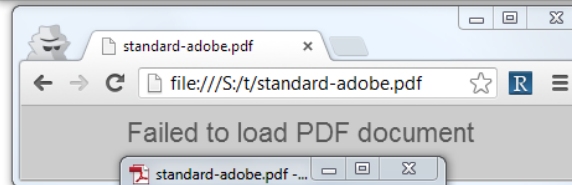
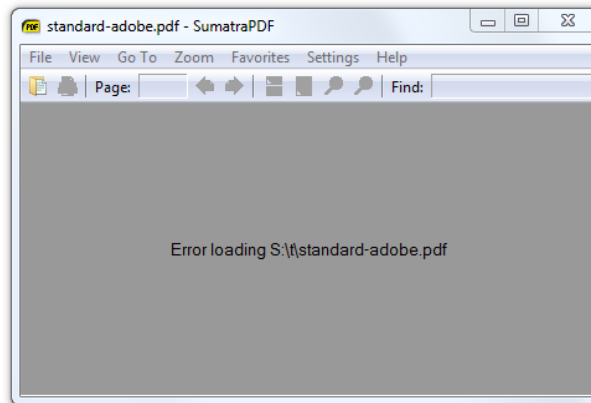
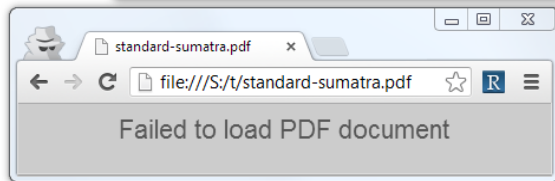
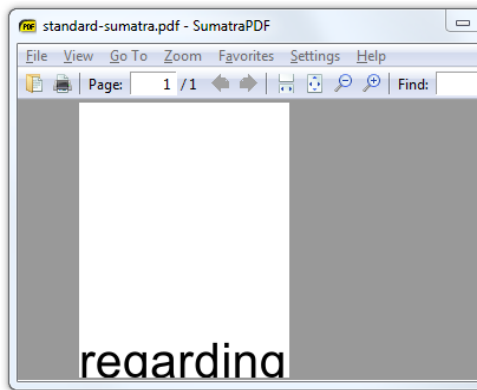


# **there is a so-called standard**

and the reality of existing parsers

looking at: Adobe, MuPDF, Chrome

- 3 different files
  - working each on a specific viewer
  - failing on the other 2



# let's look inside

- MuPDF

- no *%PDF* sig required
  - a PDF without a PDF sig ? WTF ?!?!
- no *trailer* keyword required either

- Chrome

- integer overflows:  $-4294967275 = 21$
- trailer in a comment
  - it can actually be almost ANYWHERE
  - even inside another object

- Adobe

- looks almost sane compare to the other 2

```

10 0 obj
<<

/Count 0
/Kids [<<

  /Contents 11 0 R
  /Resources <<
    /Font <<
      /F1 <<
        /BaseFont /Arial
      >>
    >>
  >>
>>]
>>

```

```

11 0 obj
<< >>
stream
BT
/F1 140
Tf
(regarding)Tj
ET
endstream

-
<</Root<</Pages 10 0 R>>>>

```

%PDF-1.

```

30 0 obj
<<

/Kids [<<
  /Parent 30 0 R
  /Contents 31 0 R
  /Resources <<>>
>>]
>>

```

```

31 0 obj
<< >>
stream
BT
/F1 150
Tf 1 0 0 1 1 0
Tm(the PDF)Tj
ET
endstream
endobj

trailer
<</Root<</Pages 30 0 R>>>>

```

%PDF

```

20 0 obj
<<
  /Pages <<
    /Kids [<<
      /Contents -4294967275 4294967296 R
    >>]
  >>
>>

```

```

21 0 obj
<< >>
stream
BT
110
Tf
('standard'...)Tj
endstream

% trailer
<</Root 20 0 R>>

```



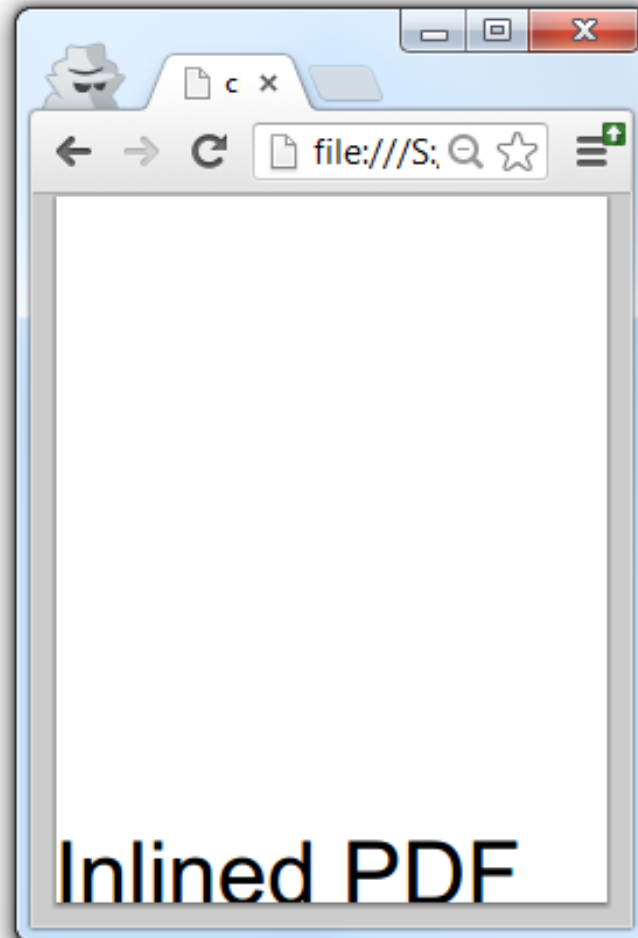
# Chrome insanity++

(thx to Jonas Magazinius)

- a single object
- no 'trailer'
- inline stream
- brackets are not even closed
- \* are required - it just checks for minimum space

All streams must be indirect objects (see Section 3.2.9, “Indirect Objects”) and the stream dictionary must be a direct object. The keyword **stream** that follows the stream dictionary should be followed by an end-of-line marker consisting of either a carriage return and a line feed or just a line feed, and not by a carriage

```
%PDF*****  
1 0 obj  
<<  
  /Size 2  
  /W[[]1/]  
  /Root 1 0 R  
  /Pages<<  
    /Kids[<<  
      /Contents<<>>  
      stream  
      BT{99  
      Tf{Td(Inlined PDF)'  
      endstream  
    >>]  
  >>  
>>  
stream  
*  
endstream  
startxref%*****
```



# PDF.JS

- very strict
  - 'too' strict / naive ?
  - I don't want to be their QA ;)
- requires a lot of information usually ignored
  - xref
  - /Length

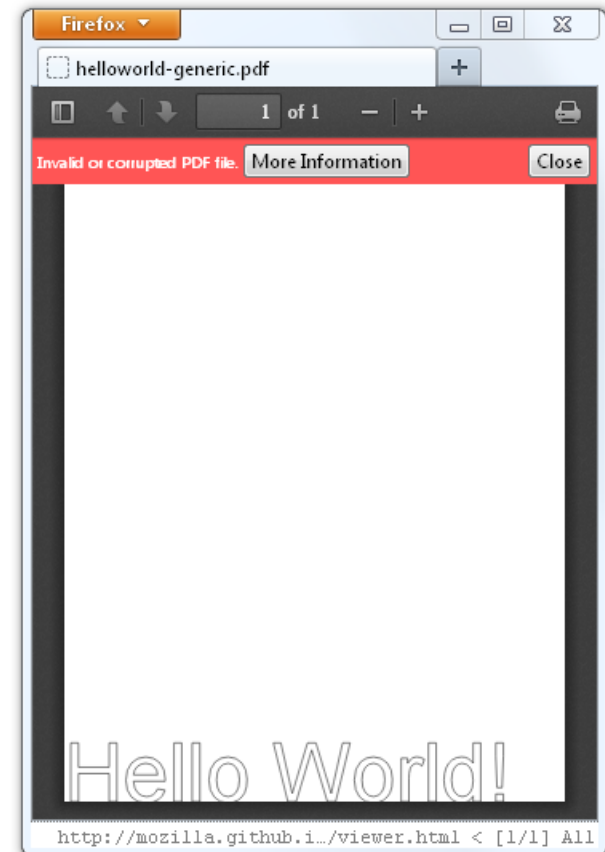
```
%PDF-1.1
1 0 obj
<<
%           /Type /Catalog
...
>>
endobj

2 0 obj
<<
           /Type /Pages
...
>>
endobj

3 0 obj
<<
           /Type /Page
           /Resources <<
             /Font <<
               /F1 <<
                 /Type /Font
                 /Subtype /Type1
...
               >>
             >>
           >>
>>
endobj

4 0 obj
<< /Length 47 >>
stream
...

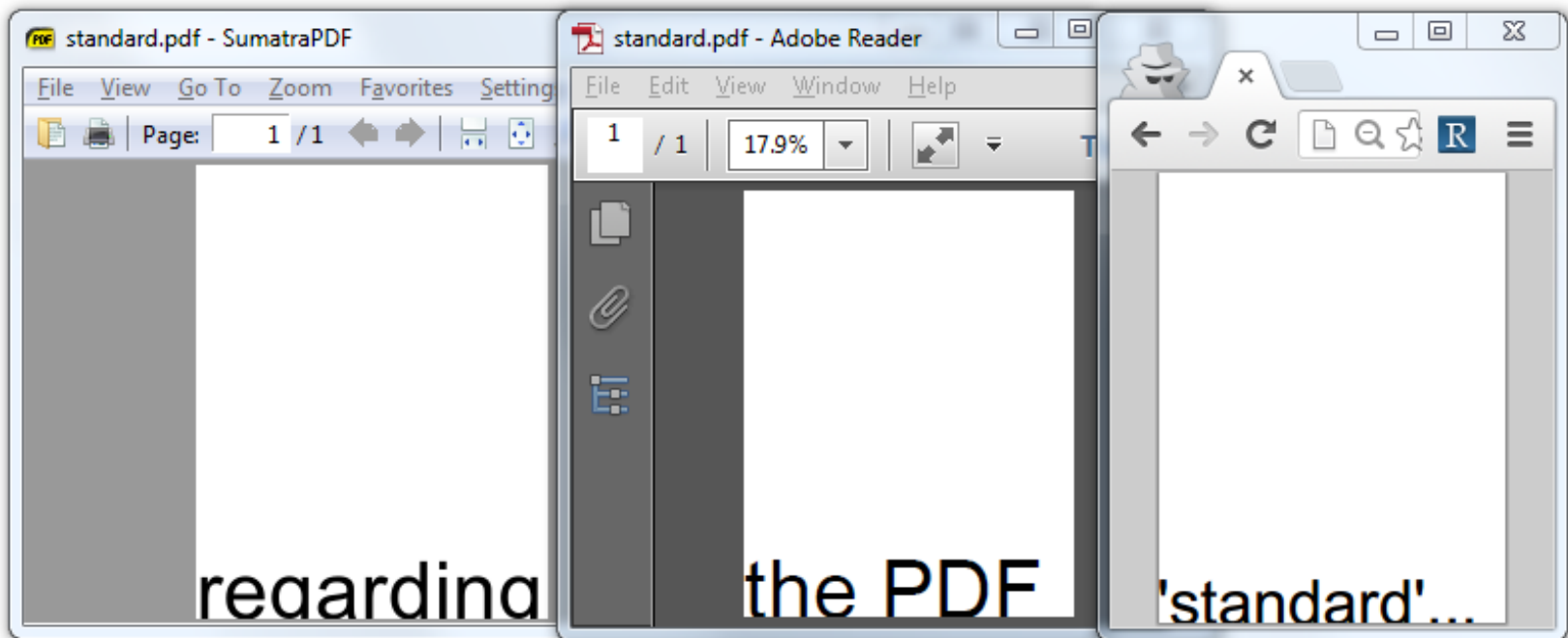
xref
0 1
0000000000 65535 f
0000000010 00000 n
...
```



# let's play further

combine 3 documents in a *single* file

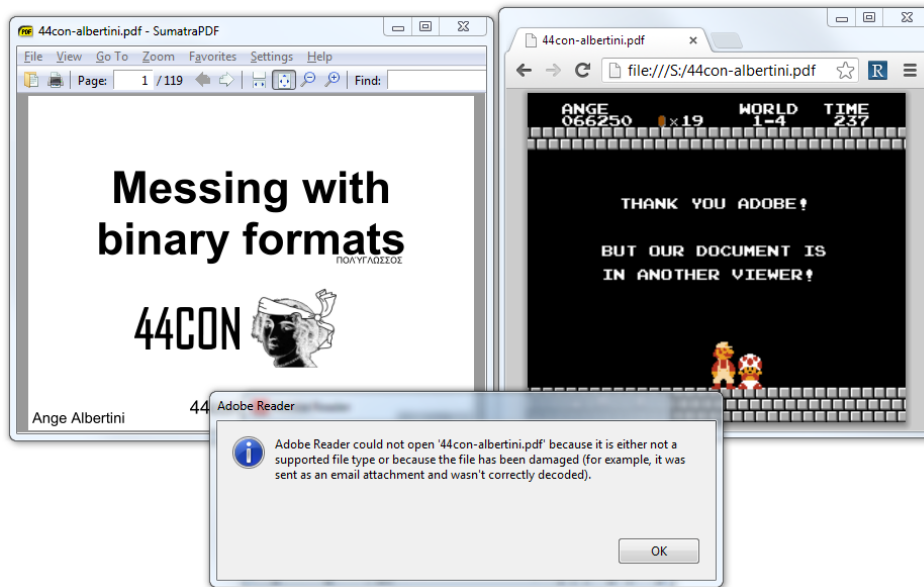
- it's actually 3 set of 'independant' objects
- objects are parsed
  - but not used



# alternate reality demo

the live slide-deck contains 2 PDF

- bogus one under Chrome
- real one under MuPDF (Sumatra, Linux...)
- rejected under Acrobat
  - because of the PE signature (see later)



alternate  
**REALITY**  
DEMO

# final PoC

- combine most previously mentioned tricks
- many fails on many tools
- total control of the structure
  - the PDF 'ends' in the Java class

```
>corkamix.exe
CorkaMIX [PE]
>java -jar corkamix.exe
CorkaMIX [java CLASS in JAR]

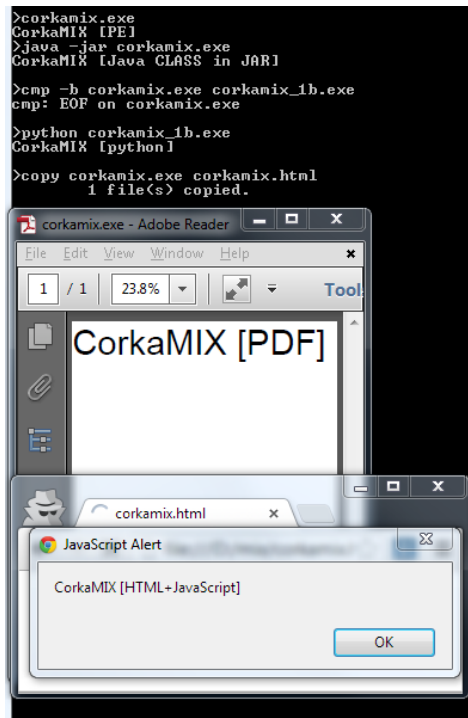
>cmp -b corkamix.exe corkamix_1b.exe
cmp: EOF on corkamix.exe

>python corkamix_1b.exe
CorkaMIX [python]

>copy corkamix.exe corkamix.html
1 file(s) copied.
```

```
db 'MZ'
; [...]
db '%PDF-1.', 0ah
db 'obj<>>stream', 0ah

db '<html>'
; [...]
    at IMAGE_NT_HEADERS.Signature, db 'PE',0,0
; [...]
db 0fh, 018h, 11b << 3
push msg
call [__imp__printf]
salc
; [...]
header:
    db 'PK', 3, 4
    dw 0ah ; version_needed
; [...]
dd 0CAFEBABEh ; signature
dw 3 ; major version
dw 2dh ; minor version
; [...]
    dd 9 ; length of bytecode
    GETSTATIC 8
    LDC 14
    INVOKEVIRTUAL 16
    RETURN
    dw 0 ; exceptions_count
    dw 0 ; attributes_count
; [...]
```



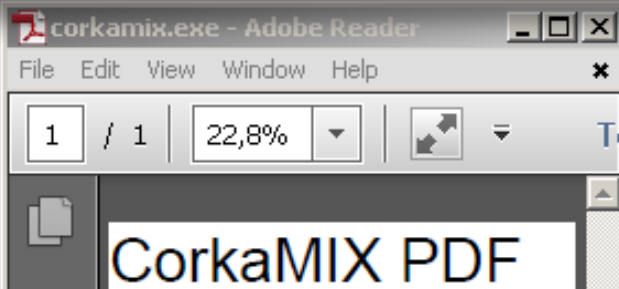
# Adobe rejects 'weird magics' after 10.1.5

not in their own specs :p

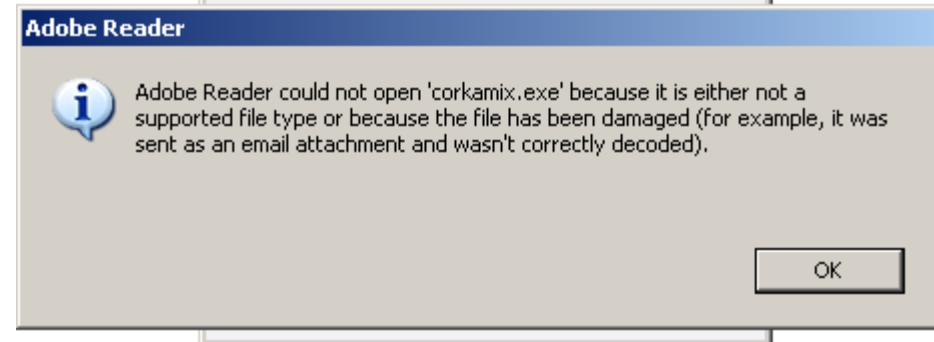
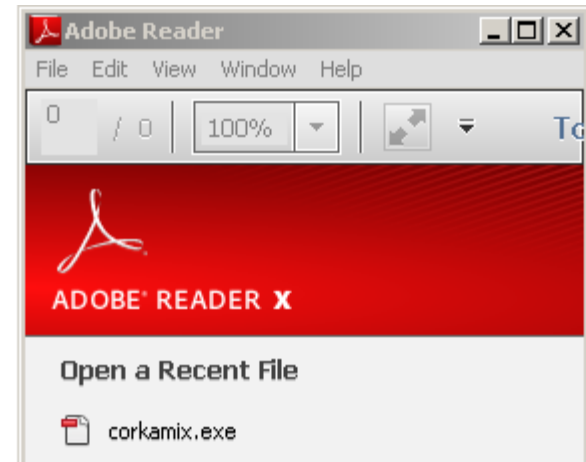
Portions copyright Right Hemisphere, Inc.

Version 10.1.4

Legal



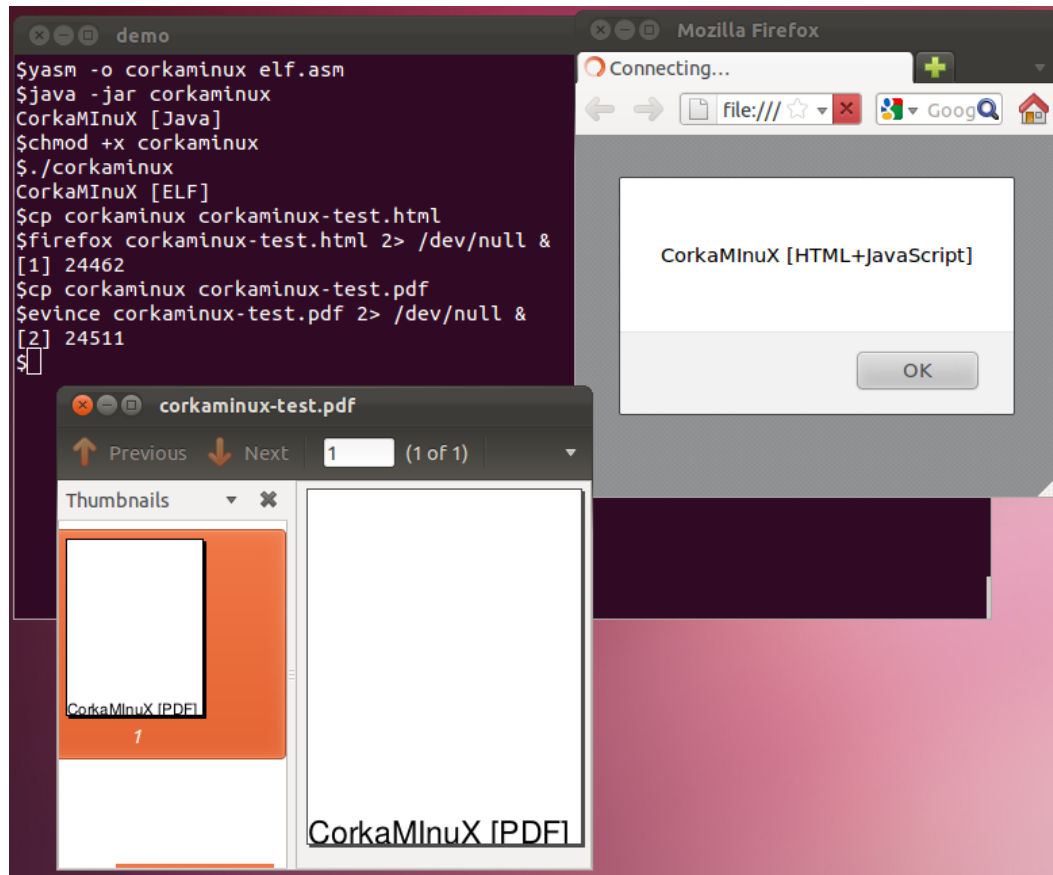
10.1.4



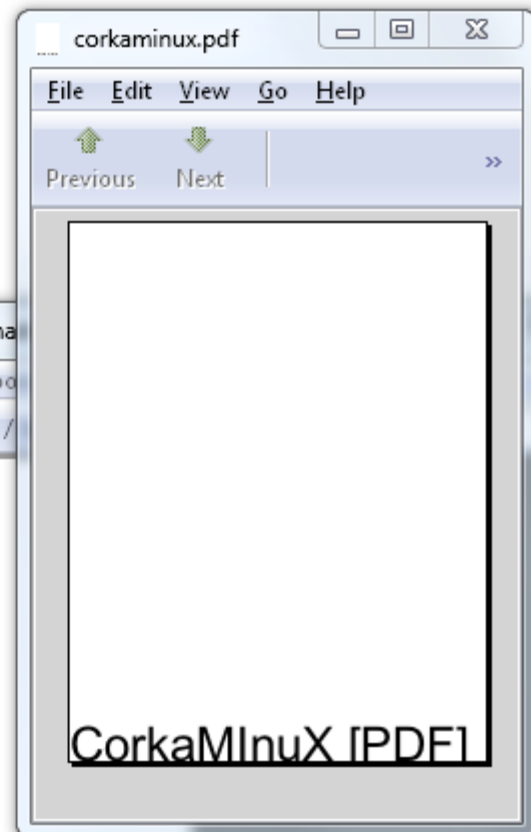
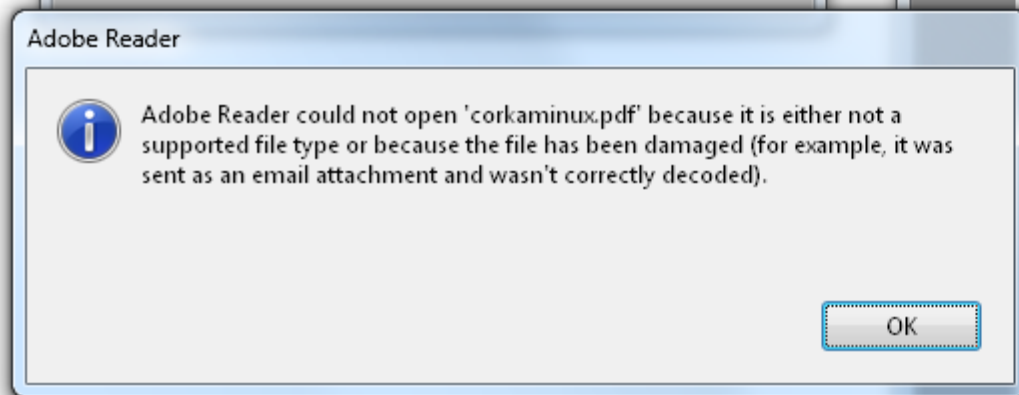
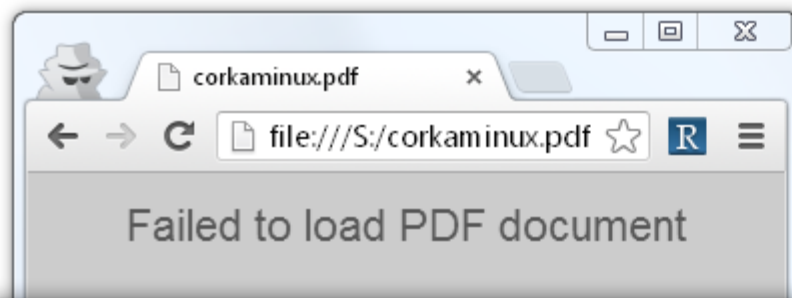
10.1.5

# also in ELF/Linux flavor

- starring a signature-less PDF
  - which won't run on other viewers

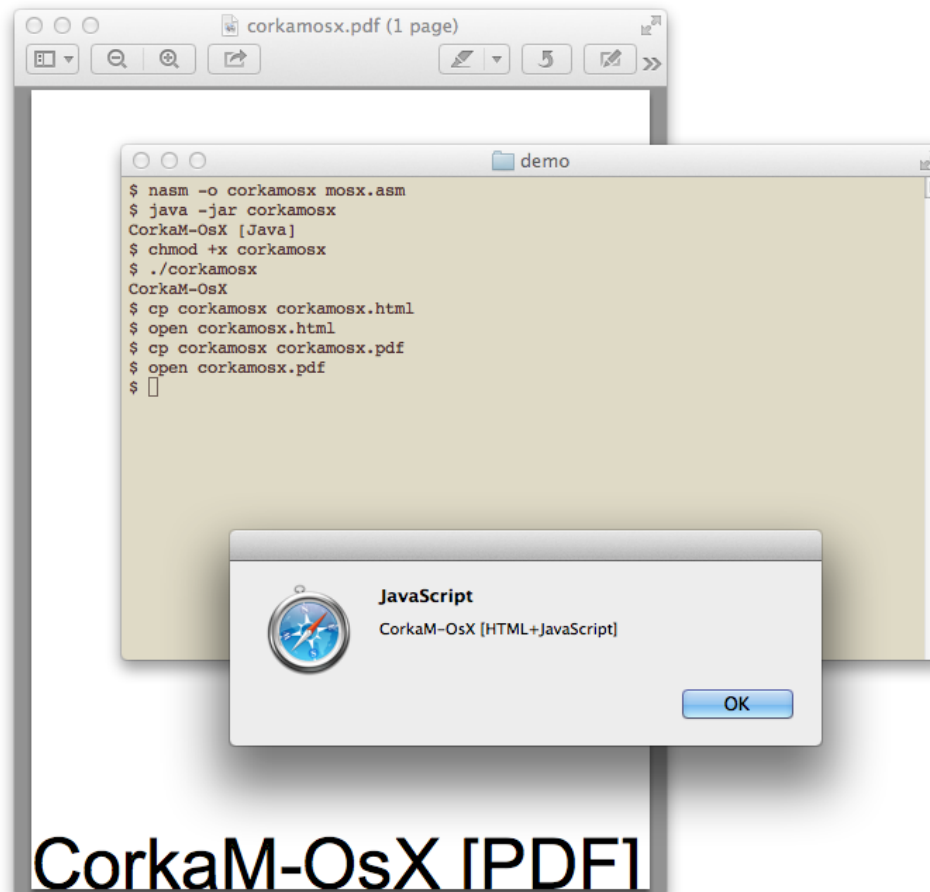






# and Apple too

PS: I don't have a Mac, this was built blindly  
Thanks to Nicolas Seriot for testing



**why should we care?**

# like washing powders

security tools are selected:

- speed
- $\{\text{files}\} \rightarrow \{[\text{clean/detected}]\}$

file types not taken into consideration



# type confusion

make the tool believe it's another type, which will fool the engine

engine with checksum caching will be fooled:

1. scanned as HTML, clean
2. reused as PE but malicious





SHA256: 2a9c7a16cdb3c3f2285afaf61072dd5e7cc022e97f351cad6234a13e5216f389

SHA1: e27faaa006229f8e4ab97fba7019dc9f2797f84d

MD5: 88cad2b56ab67b43794a0f7a4e690fd5

File size: 1.5 KB ( 1530 bytes )

File name: corkamix.exe

File type: PDF

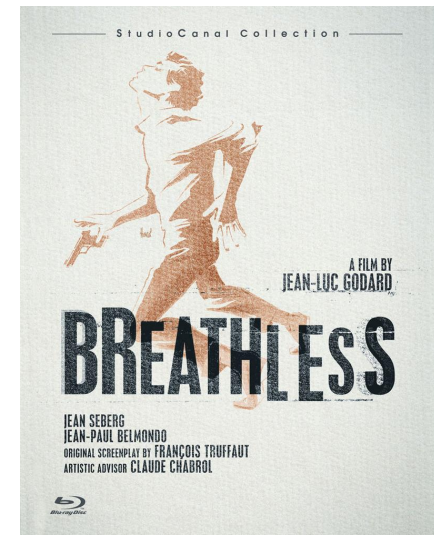
Tags: pdf

# engine exhaustion

rankings in magazines are based on scanning time

→ scanning per file must stop arbitrarily

→ waste scanning cycle by adding extra formats



# Weaknesses

- evasion
  - filters → exfiltration
  - *same origin policy*
  - detection
    - ex: clean PE but malicious PDF/HTML/...
    - exhaust checks
    - pretend to be corrupt
- DoS



# Conclusion

# Conclusion

- type confusion is bad
  - succinct docs too
  - lazy softwares as well
- go beyond the specs
  - Adobe: good
- suggestions
  - more extensions checks
  - isolate downloaded files
  - enforce magic signature at offset 0

thank YOU !

**Questions ?**

http://

reverseengineering

.stackexchange.com

@angealbertini



ange@corkami.com

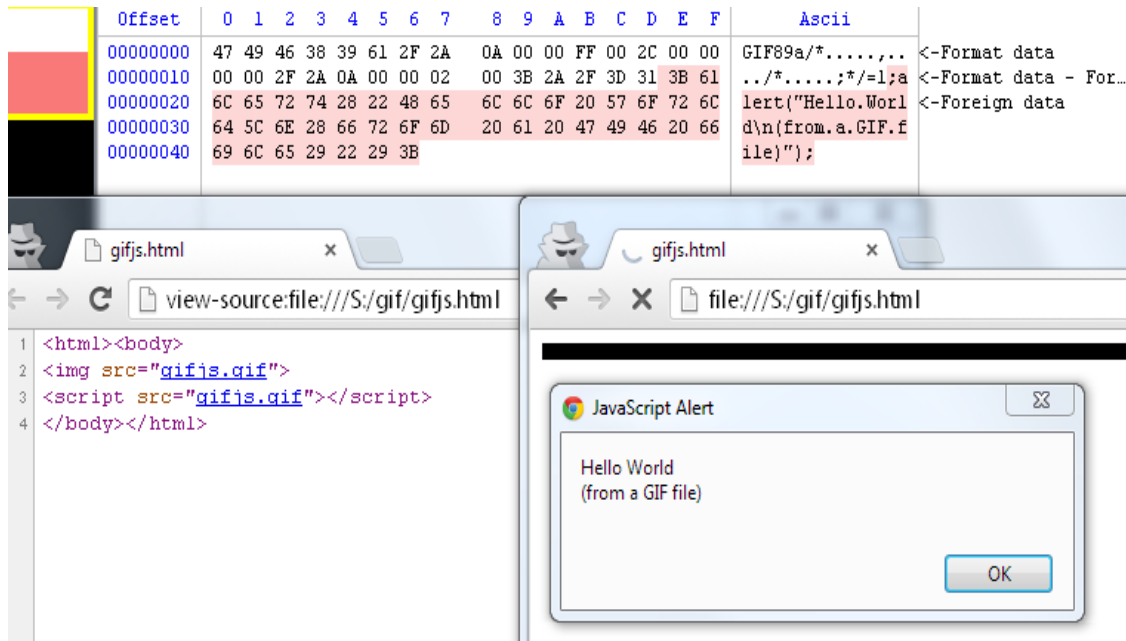
**Bonus**

# Valid image as JavaScript

Highlighted by Saumil Shah

- abusing header and parsers laxisms
- turn a field into /\*
- close comment after the picture data

Offset	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F	Ascii
00000000	47	49	46	38	39	61	2F	2A	0A	00	00	FF	00	2C	00	00	GIF89a/*.....,...
00000010	00	00	2F	2A	0A	00	00	02	00	3B	2A	2F	3D	31	3B	61	../*.....;*/=1;a
00000020	6C	65	72	74	28	22	48	65	6C	6C	6F	20	57	6F	72	6C	lert("Hello.Worl
00000030	64	5C	6E	28	66	72	6F	6D	20	61	20	47	49	46	20	66	d\n(from.a.GIF,f
00000040	69	6C	65	29	22	29	3B										ile)");

```
1 <html><body>
2 
3 <script src="gifjs.gif"></script>
4 </body></html>
```