

# Formaty plików wykonywalnych dla systemu DOS, Windows i Linux

Mateusz Gorządek i Arkadiusz Grzanka



# Formaty plików wykonywalnych dla systemu DOS, Windows i Linux

Mateusz Gorządek i Arkadiusz Grzanka



# Plik wykonywalny

Plik, który może być uruchomiony bezpośrednio w środowisku systemu operacyjnego.

Najczęściej zawiera binarną reprezentację instrukcji konkretnego typu procesora. Oprócz tego znajdują się w nim zwykle wywołania systemowe, dlatego pliki wykonywalne zazwyczaj są specyficzne nie tylko dla danego procesora, ale też dla danego systemu operacyjnego. Nie dotyczy to sytuacji, kiedy zawierają formę pośrednią, która do uruchomienia wymaga interpretera lub maszyny wirtualnej – takie pliki mogą być zwykle uruchamiane na różnych systemach.

Specyficzną odmianą plików wykonywalnych są skrypty powłoki. Zarówno w systemie DOS jak i systemach uniksowych pliki takie można uruchamiać bezpośrednio, jak każdy inny program (np. wpisując jego nazwę w wierszu poleceń).

Zależnie od konwencji, nazwy plików wykonywalnych mogą wyróżniać się rozszerzeniem, np. w DOS i Windows przyjęte zostało rozszerzenie COM i EXE (stąd popularne określenie "egzek"). W systemach uniksowych pliki mają ustawiony atrybut wykonywalności (oznaczany literą x).

# COM (COre iMage)

Struktura plików COM jest bardzo prosta, nie zawiera nagłówka ani metadanych, tylko sam kod i dane. Plik binarny ma maksymalnie 65280 (0xFF00 w systemie szesnastkowym) bajtów (nieco mniej niż 64 kB), wszystkie dane oraz kod są trzymane w jednym segmencie. Jest to dość poważną wadą, ponieważ w takim przypadku programy COM nie mogą być łączone z modułami zewnętrznymi.

Programy tego typu są ładowane przez system operacyjny zawsze pod stały adres offsetowy 0x0100, skąd są wykonywane. Nie było to problemem w przypadku ośmiobitowych maszyn, ale jest to głównym powodem dlaczego ten format wyszedł z użycia zaraz po wprowadzeniu 16-bitowych i 32-bitowych procesorów z większymi, segmentowanymi pamięciami, gdzie każdy proces dostaje swój własny segment.

Lista parametrów uruchamianego programu umieszczana jest z przesunięciem 0x0081 (ze spacją), poprzedza ją 8 bitowa długość w offsecie 0x0080 w tym samym segmencie.

# MZ DOS

16-bitowy format plików wykonywalnych .EXE w systemie DOS.

Identyfikowany za pomocą dwóch znaków „MZ” (hexadecimal: 4D 5A) ASCII na początku pliku. „MZ” to inicjały nazwiska Marka Zbikowskiego, jednego z programistów systemu MS DOS. Format MZ DOS jest nowszy od COM (zastąpił COM w MS-DOS 2.0) i posiada od niego nieco większe możliwości.

Nagłówek formatu DOS MZ zawiera informacje na temat położenia poszczególnych segmentów, co pozwala na rozmieszczenie ich na dowolnych adresach w pamięci. Wspierane są również pliki wykonywalne większe niż 64 KiB.

Egzeki MZ DOS mogą być uruchamiane w DOS oraz systemach operacyjnych z rodziny Windows 9x. Na 32-bitowych systemach serii Windows NT możliwe jest ich uruchomienie za pomocą wbudowanej maszyny wirtualnej Virtual DOS. Nie można ich uruchomić na 64-bitowej wersji systemu Windows.

# NS (New Executable)

16-bitowym format plików .EXE i następca formatu DOS MZ.

Używany był w systemie DOS 4.0 i w 16-bitowych wersjach Microsoft Windows (Windows 1.0, Windows 2.x i Windows 3.x).

Format ten zawiera nagłówek pliku DOS MZ (tzw. DOS stub) w celach wstecznej kompatybilności z nieobsługiwanymi wersjami DOSa.

Pliki wykonywalne formatu NE mogą być uruchamiane na 32-bitowych Windowsach. Nie są natomiast wspierane na Windowsach 64-bitowych.

# PE (Portable Executable)

Format plików wykonywalnych, obiektowych oraz bibliotek dynamicznych. Używany w 32- i 64-bitowych wersjach systemów operacyjnych z rodziny Microsoft Windows. Portable oznacza "przenośny", co odnosi się do uniwersalności formatu, dostępnego w wielu architekturach systemów komputerowych.

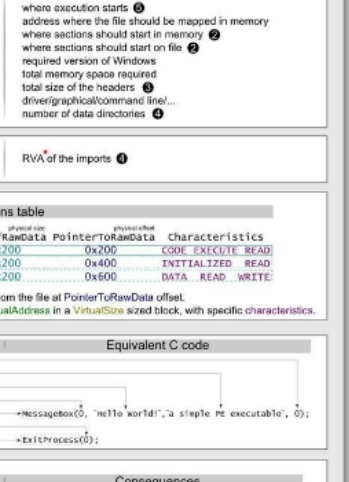
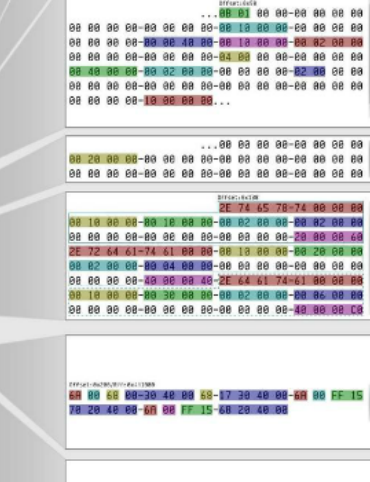
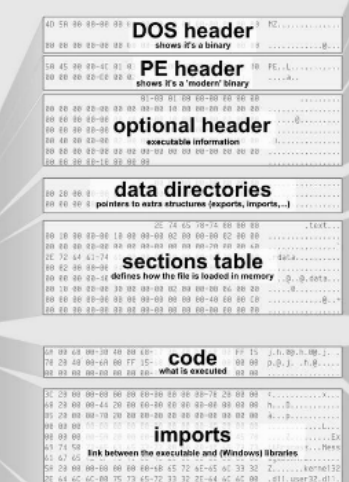
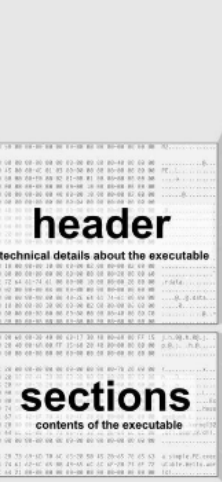
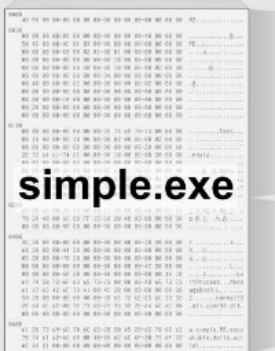
W systemach z rodziny Windows NT, format PE jest używany m.in. przez pliki \*.exe, \*.dll, \*.obj, \*.sys (najczęściej plik sterownika urządzenia).

Format PE jest zmodyfikowaną wersją Uniksowego formatu COFF stąd też jego alternatywna nazwa – PE/COFF.

W systemach Windows NT, format PE może zawierać zarówno instrukcje z zestawu IA-32 jak i IA-64 oraz x86\_64 (AMD64 i EM64T). Do wersji 4 włącznie, Windows NT (a więc de facto PE) obsługiwał również architektury MIPS, DEC Alpha i PowerPC. PE używany jest również w Microsoft Windows CE, który kontynuuje wsparcie dla kilku wariantów architektury MIPS, ARM (włączając Thumb) oraz SuperH.



## Dissected PE



Hexadecimal dump	ASCII dump	Fields	Values	Explanation
4D 5B 00	MZ..... .....0...	e_magic e_lfanew	'MZ' 0x40	constant signature offset of the PE Header ❶
5A 45 00	PE..L..... .....	Signature Machine NumberOfSections SizeOfOptionalHeader Characteristics	'PE', 0, 0 0x14c [intel 386] 3 0x0 0x102 [32b EXE]	constant signature processor: ARM/MIPS/intel... number of sections ❷ relative offset of the section table ❷ EXE/DLL...
00 00	..... ..... ..... ..... ..... .....	Magic AddressOfEntryPoint ImageBase SectionAlignment FileAlignment MajorSubsystemVersion SizeOfImage SizeOfHeaders Subsystem NumberOfRvaAndSizes	0x10b [32b] 0x1000 0x400000 0x1000 0x200 4 [NT 4 or later] 0x4000 0x200 2 [GUI] 16	32 bits/64 bits where execution starts ❸ address where the file should be mapped in memory where sections should start in memory ❹ where sections should start on file ❹ required version of Windows total memory space required total size of the headers ❺ driver\graphical\command line... number of data directories ❻
00 00	..... .....	ImportsVA	0x2000	RVA of the imports ❶

Sections table						
Name	VirtualSize	VirtualAddress	SizeOfRawData	PointerToRawData	Characteristics	
.text	0x1000	0x1000	0x200	0x200	CODE_EXECUTE_READ	
.pdata	0x1000	0x2000	0x200	0x400	INITIALIZED_READ	
.data	0x1000	0x3000	0x200	0x600	DATA_READ_WRITE	

For each section, a **SizeOfRawData** sized block is read from the file at **PointerToRawData** offset. It will be loaded in memory at address **ImageBase + VirtualAddress** in a **VirtualSize** sized block, with specific characteristics.

x86 assembly	Equivalent C code
<pre> push 0 push 0x403000 push 0x403017 push 0 call [0x402070] push 0 call [0x402065]                     </pre>	<pre> MessageBox(0, "hello world!", "a simple PE executable", 0);                     </pre>

Imports structures	Consequences
<pre> descriptors 0x203C -&gt; 0x204c, 0 Hint 0x2078 -&gt; kernel32.dll, 0 Hint, 0 ExitProcess 0x2068 -&gt; 0x204c, 0 Hint, 0 ExitProcess 0x2044 -&gt; 0x205a, 0 Hint, 0 MessageBoxA 0x2085 -&gt; user32.dll, 0 Hint, 0 MessageBoxA 0x2070 -&gt; 0x205a, 0 Hint, 0 MessageBoxA                     </pre>	<p>after loading, 0x1c2068 will point to kernel32.dll's ExitProcess 0x1c2070 will point to user32.dll's MessageBoxA</p>

Strings
<pre> a simple PE executable\0 hello world!\0                     </pre>

**simple.exe**  
technical details about the executable

**header**  
technical details about the executable

**sections table**  
defines how the file is loaded in memory

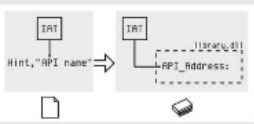
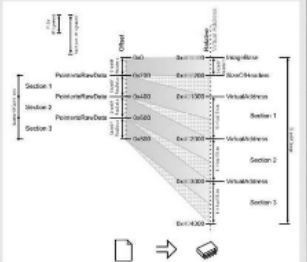
**code**  
what is executed

**imports**  
link between the executable and (Windows) libraries

**data**  
information used by the code

## Loading process

- 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 offset (OptionalHeader) + SizeOfOptionalHeader)  
it contains NumberOfSections elements  
it is checked for validity with alignments:  
Alignments 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 NumberOfRvaAndSizes  
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



## Notes

- MZ HEADER** aka DOS\_HEADER  
Starts with "MZ" (initials of Mark Zbikowski MS-DOS developer)
- PE HEADER** aka IMAGE\_FILE\_HEADERS / COFF file header  
Starts with "PE" (Portable Executable)
- OPTIONAL HEADER** aka IMAGE\_OPTIONAL\_HEADER  
Optional only for non-standard PE's but required for executables
- RVA** Relative Virtual Address  
Address relative to ImageBase (at ImageBase, RVA = 0)  
Almost all addresses of the headers are RVAs  
In code, addresses are not relative.
- INT** Import Name Table  
Null-terminated list of pointers to Hint, Name structures
- IAT** Import Address Table  
Null-terminated list of pointers  
On file it is a copy of the INT  
After loading it points to the imported APIs
- HINT**  
Index in the exports table of a DLL to be imported  
Not required but provides a speed-up by reducing look-up



# Pliki wykonywalne w systemie Linux

Linux wspiera kilka różnych formatów plików wykonywalnych. Każdy z nich zarejestrowany jest na liście formats.

```
struct linux_binfmt {
    struct linux_binfmt * next;
    long *use_count;
    int (*load_binary)(struct linux_binprm *, struct pt_regs * regs);
    int (*load_shlib)(int fd);
    int (*core_dump)(long signr, struct pt_regs * regs);
};
```

Każdy z formatów za pomocą funkcji `register_binfmt` definiuje i wstawia do listy `formats` elementy typu `linux_binfmt` z odpowiednimi dla siebie funkcjami do rozpoznawania i ładowania programu, ładowania bibliotek dzielonych i zrzutu obrazu procesu na dysk.

4 rodzaje formatów plików wykonywalnych: skrypty, format Javy, a.out, ELF

# Formaty skryptów

Wykonywane są pliki tekstowe napisane w dowolnym języku skryptowym, jakiego interpreter znajduje się w systemie.

Plik musi mieć ustawiony atrybut do wykonania. Można tego dokonać za pomocą polecenia:

```
chmod +x nazwa_pliku
```

Ścieżka interpretera zdefiniowana jest w pierwszej linii skryptu, zaraz po znakach `#!`. Np. `#!/bin/bash/`

# Formaty Javy

W Linuxie istnieje możliwość bezpośredniego wywoływania plików napisanych w javie.

Spełnione muszą zostać następujące warunki:

- w systemie musi być zainstalowany Java Developers Kit for Linux
- samodzielne programy-klasy, które mają być bezpośrednio wykonywane muszą mieć:
  - ustawiony atrybut do wykonania
  - zaimplementowaną metodę `public void main(String args[])`.
- pliki, które zawierają aplety i mają być bezpośrednio wykonywane muszą mieć:
  - następującą pierwszą linię w pliku: `<!--applet-->`,
  - ustawiony atrybut do wykonania.
- jeśli klasy, których używa program-klasa bądź aplet, znajdują się poza katalogiem `/usr/local/java/classes/` , to należy eksportować zmienną środowiskową `CLASSPATH` z odpowiednio zmodyfikowanymi parametrami.

# a.out

Format plików używany szczególnie w Uniksie i innych uniksopodobnych systemach operacyjnych jako format plików wykonywalnych, plików obiektowych oraz bibliotek dzielonych.

a.out jest wciąż domyślną nazwą pliku wynikowego w niektórych kompilatorach

Był to podstawowy format plików wykonywalnych w Linuxie.

Obsługa tego formatu przez jądro zapewnia:

- ładowanie programu
- ładowanie dzielonej biblioteki procedur
- zrzut obrazu na dysk.

# a.out

Plik wykonywalny w formacie a.out zawiera cztery części:

- Nagłówek podstawowy opisujący ilość sekcji w pliku
- Nagłówki sekcji
- Sekcje z kodem
- Sekcje z tablicą symboli i innymi danymi użytecznymi podczas uruchamiania programu.

struct exec
text
data
(bss)
(symbol table etc.)

a.out file format

# ELF

ELF - Executable and Linking Format

Format plików binarnych stworzony przez USL (UNIX System Laboratories)

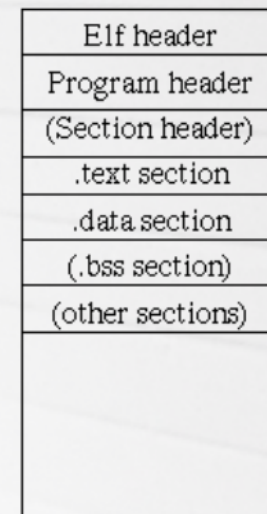
Posiada następujące możliwości:

- plik wykonywalny może być podzielony na dowolną liczbę sekcji
- sekcje mają oddzielnie określone cechy
- mogą być zaznaczone jako wcale nieładowne do pamięci
- możliwość przeszukiwania funkcji w bibliotekach po kompilacji
- duża elastyczność i modyfikowalność
- tworzenie dzielonych bibliotek jest dużo łatwiejsze
- łatwość dynamicznego wczytywania (np. dla programów które wczytują dodatkowe moduły w trakcie działania).

# ELF

Plik ELF składa się z:

- jednego nagłówka programu
- listy nagłówków programu zawierającej zero lub więcej segmentów
- listy nagłówków sekcji zawierającej zero lub więcej sekcji
- danych zawierających segmenty i sekcje



ELF file format

## Porównanie formatów plików wykonywalnych

↕	Explicit processor declarations	Arbitrary sections	Metadata <sup>[vague]</sup>	Digital signature	String table	Symbol table	64-bit	Fat binaries	Can contain icon
a.out	No	No	No	No	Yes <sup>[1]</sup>	Yes <sup>[1]</sup>	Extension	No	No
COFF	Yes by file	Yes	No	No	Yes	Yes	Extension	No	No
ELF	Yes by file	Yes	Yes	Yes <sup>[2]</sup>	Yes	Yes <sup>[3]</sup>	Yes	Extension <sup>[4]</sup>	Extension <sup>[5]</sup>
PE	Yes by file	Yes	Yes	Yes <sup>[6]</sup>	Yes	Yes	Yes	No	Yes
Mach-O <sup>[7]</sup>	Yes by section	Some (limited to max. 256 sections)	Yes	Yes	Yes	Yes	Yes	Yes	No
SOM	Unknown	Unknown	No	No	Unknown	Yes	No	Unknown	No
Hunk	Unknown	Yes	Yes	No	No	Yes	No	Yes	No
MZ	No (x86 only)	Yes	No	No	No	No	No	No	No
DOS COM	No (x86 only)	No	No	No	No	No	No	Extension	No
CP/M-80 COM	No (8080/Z80 only)	Extension (CP/M 3 only)	No	No	No	No	No	Extension	No
CP/M-86 CMD	No (x86 only)	Yes	No	No	No	No	No	No	No
PEF <sup>[8]</sup>	Yes by file	No	No	No	Yes	Yes	No	No	No
ECOFF	Yes by file	Yes	No	No	Yes	Yes	Yes	No	No
XCOFF	Yes by file	Yes	No	No	Yes	Yes <sup>[9]</sup>	Yes	No	No
NE	Unknown	Unknown	Unknown	No	Unknown	Unknown	No	No	Yes
LX	Unknown	Unknown	Unknown	Unknown	No	Yes <sup>[10]</sup>	No	No	Yes



# Formaty plików wykonywalnych dla systemu DOS, Windows i Linux

Mateusz Gorządek i Arkadiusz Grzanka

