

Metadata Since it's present in the file but not needed for parsing or rendering the file's contents, metadata is a great source of abuse. Many old school formats have fixed length fields, as hard coding was the norm back then.

Comments are typically ignored empty space with typically a set length that is declared before, and are present in most file formats. Unlike PDF, XML enforces an encoding for its comments like the rest of the file, but in general, comments are just ignored, and preserved, no matter their amount, their length or their content.

Comments are not the only source of abuse. Extensible metadata with a user-chosen ID, or fields like file names in an archive can also be used to store some foreign data. A notable exception is that in Gzip, the optional comment and file name are null-terminated, which shows that they're intended to store standard text, while the also optional Extra Field is defined with a 2-byte length.

As a side note, metadata may seem like a permanent risk entirely, and it's natural to wonder why we define them officially in every format if they are so easily abused later. While metadata doesn't seem like an initial requirement to keep the format simple — like the Quite OK Image format — it is eventually needed to be able to keep extra information in the file, which is exactly what happened for the MP3 files.

At the release of *l3enc* (the original Mpeg Layer 3 encoder in 1994), the files initially had an 13 extension, had no file format whatsoever. They were pure sequences of layer 3 frames, each with their own frame header with no signature, making them hard to detect and easy to confuse with other data such as JPEG segments.

Since there was no way to store any metadata in L3 files, the compatible ID3v1 *footer* with a hard coded length was unofficially defined. More structures were defined in other clumsy ways around the L3 stream (Xing, Lame, APE, ...), showing the need for proper definition of metadata storage from the beginning. ID3v2 eventually defined a header, a magic which gave at last a proper format to MP3 streams.

It's a shortsighted move to come up with a great compression algorithm (e.g., MP3, QOI) and define a way to store some data in a file format without the ability to extend it with new but optional data in the future. Even if it means that these structures can be abused, you can't have an extensible format

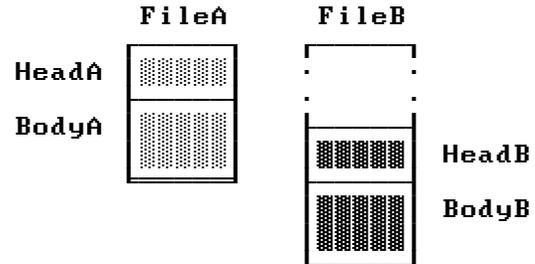
that can't be abused, and people *will* extend or fork your format if not, requiring an extra format that could have been avoided from the beginning.

Wrapping Some formats don't tolerate appended data, but they can end with a parasite-hosting structure, which acts like appended data, just wrapped in a declarative structure. Since most data-storing structures have all their declaration before their data itself, it behaves like appended data from the outside, even if the length of the appended data has to be declared somewhere. It looks like appended data — ignored but tolerated but it's technically an appended parasite — declared and skipped.

Zipper Some formats have very strong constraints: a tiny cavity, or very small parasite length (256 bytes for a GIF). Rather than parasitizing a whole file, let's just add the declaration of the file, a declaration of a comment, and then store the rest of the file as appended data.

A zipper is a fabric construct of two sliders where each tooth interlocks with the other side's tooth. As an analogy, a *zipper* is a file construct where each format comments the other format's elements, and each tooth is a parasite for the other format.

The simplest form of a zipper is made of two formats: Format A starts at offset zero, tolerates appended or wrapped data, and Format B starts with a cavity. Both formats can be parasitized.



SWTP 6800 OWNERS—WE HAVE A CASSETTE I/O FOR YOU!

The CIS-30+ allows you to record and playback data using an ordinary cassette recorder at 30, 60 or 120 Bytes/Sec. No Hassle! Your terminal connects to the CIS-30+ which plugs into either the Control (MP-C) or Serial (MP-S) Interface of your SWTP 6800 Computer. The CIS-30+ uses the self clocking 'Kansas City' Biphase Standard. The CIS-30+ is the FASTEST, MOST RELIABLE CASSETTE I/O you can buy for your SWTP 6800 Computer.

PerCom has a Cassette I/O for your computer!
Call or Write for complete specifications



Kit — \$69.95*
Assembled — \$89.95*
(manual included)
* plus 5% t/shipping



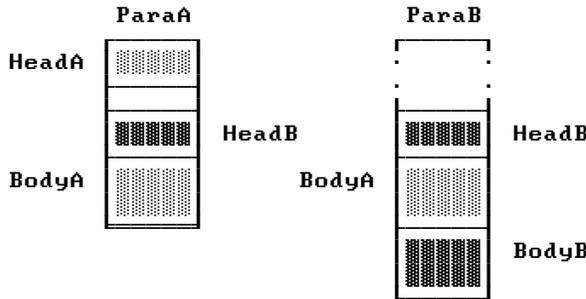
PerCom Data Co.
P.O. Box 40598 • Garland, Texas 75042 • (214) 276-1968

PerCom — "peripherals for personal computing"

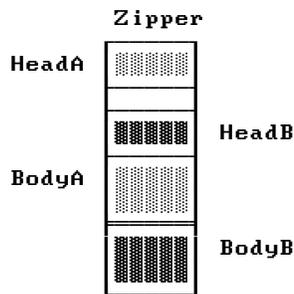


TEXAS RESIDENTS ADD SALES TAX

We parasitize File A with Head B, adding padding if required. We also parasitize File B with Body A, *wrapping* Body B in advance if required.



When we merge these files, it looks like this.



Zipper combines various format features (cavity, parasite, appended data) to overcome limitations and make even more weird formats combine.

Mitra is a tool that combines all this knowledge for 40+ different format, generating hundreds of format combinations with different strategies.⁷

Mitra is a simple tool. It doesn't understand file formats structure, it just contains the minimum amount of information to identify and parasitize a file format. It expects standard files as input!

Abuses

Payload embedding The simplest form of exploitation is to just embed a payload that doesn't need to be a valid file. In this case, use the `--force` command line parameter.

The universal example for that is HTML or JavaScript that can be embedded in most file formats. If the file is too big, the HTML page might take too long to load entirely. In that case, use JavaScript to break out of the appended data and limit the parsing to the web payload only.

⁷git clone https://github.com/corkami/mitra

⁸unzip pocorgtfo22.pdf mocky.py

Mocks While 7:6 covered file type identification, it didn't cover any exploitation. The easiest way to exploit file type identification is just to give a binary blob the right signature at the right offset. It could even happen accidentally. Such a file is a *mock* file. A simple example is `FF D8`, a two byte file that can be identified as a JPEG image.

Polymocks It's also interesting to just add a mock signature at a given offset in a valid file via any of the previously mentioned techniques.

Mocky is a tool that uses the Mitra library to insert specific filetype signatures at specific offsets to create polymocks.⁸ Since programs like `file` have so many signatures that they are scanned by alphabetic order of their category, it's possible to predict which detection will be returned first, and the order might not be what you'd expect from a threat model perspective.

Of course, it's possible to cram as many signatures as possible within the constraint of the target format, such as this issue, a valid PDF with many extra mock signatures stored in a standard stream object. Here is an example of such a *polymock* file, with many filetype detections yet no valid content:

	x0 x1 x2 x3 x4 x5 x6 x7 x8 x9 xa xb xc xd xe xf	0+2 Dos executable
0x	M Z 50 EA j P 01 07 19 04 00 10 5 N D H	2+2 Arj
		4+2 JPEG 2000
		6+2 UnixOS
		8+4 Symbian
		C+4 Sndh
1x	N R 0 0 DC A7 C4 FD 5D 1C 9E A3 R E - A	10+4 Nintendo Switch
		14+4 Zoo
		18+4 Nintendo Wii
		1C+4 Rar v3.4
		20+4 AFS
2x	N X 5 B 18 28 6F 01 P K 03 04 P T M F	24+4 zImage
		28+4 PKZip
3x	S Y M E x e 7 z BC AF 27 1C 5 0 N G	2C+4 PolyTracker
		30+6 SymbOs
		36+6 7-zip
		3C+4 SoundFX
4x	7F 10 DA BE 00 00 CD 21 P K 01 02 5 C R 5	40+4 VirtBox
		46+2 Int 21h
		48+4 PKZip
		4C+4 ScreamTracker
5x	R a F ! 1A 07 01 \0 L R Z I P L O T	50+8 Rar v5
6x	X Z 8 4	58+4 LzZip
		5C+8 Plot84
x4	...	64+7 Rar v4
7x	...	6D+5 EZD Map
x2	...	72+6 Xz
		78+4 LZ4
		7C+4 LZ4
8x	D I C M % P D F - 1 . 4 \n o b j	80+4 DBCOM
		84+C PDF

Running `file` with `--keepreading` gives an impressive list of detected formats:

```
Plot84 plotting file
SymbOS executable v7.z, name: ...
Old EZD Electron Density Map
Zoo archive data, vj., modify: v78.88+
Symbian installation file
Scream Tracker Sample adlib drum stereo...
Poly Tracker PTM Module Title: "MZ..."
SoundFX Module sound file
Nintendo Wii disc image: "NXSB..."
DICOM medical imaging data
Linux kernel ARM boot executable...
VirtualBox Disk Image, minor 8653 (MZ...)
JPEG 2000 image
ARJ archive data
COM executable for DOS
unicos (cray) executable
data
```

Mocky has a `--combine` flag for to try and insert as many signatures in a file as possible.

`file` has an extra weakness that it has special support for tar files before any other format, and can identify tar files not by their magic, but by the validity of their header checksum, even without any tar signature in the file. This is used in this issue too, so even if the file is a standard PDF starting with a generic PDF signature, `file` with no parameter sees it as a `tar archive`, even if it doesn't contain a magic Tar signature.

Mocky will adjust the tar checksum if used for a polymock file. Here is such an empty `mock.tar` file, detected generically as `tar archive` even if it contains no signature at all:

```
000: 00000000 00000000 00000000 00000000
...
090: 00000000 4 0 0      00  00000000
...
1F0: 00000000 00000000 00000000 00000000
```

Adding a valid tar checksum to the previous polymock example will indeed return a tar filetype — if the `--keep reading` parameter isn't used — despite all the other present signatures.

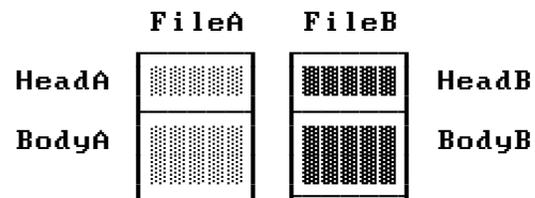
Near-polyglots Formats that require a different signature at the same offset can't be combined in a polyglot. However their combination can still be exploited in different conditions. *Near polyglots* are files that are almost polyglots, except that some bytes have to be replaced so that the file type changes.

This change could happen over the network if some packets arrive in a different order. It could also happen due to weak bits, leading to different contents. And it can happen via a cryptographic

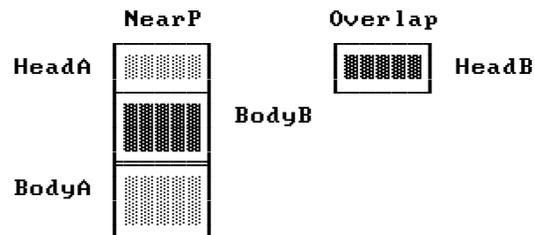
operation in which case you can call them *crypto-polyglots*.

One of these use case is Angecrption, introduced in PoC||GTFO 3:11, where I demonstrated abusing the Initialization Vector of CBC, CFB or OFB block modes to replace the first block of the crypto-polyglot. A new abuse of pseudo-polyglots is presented in the TimeCrption article on page 30 of this release.

To generate a near polyglot, you need very light constraints. `FormatA` can be parasitized and `FormatB` can start at the same offset, and needs to tolerate appended or wrapped data. Technically, generating a near polyglot is like parasitizing a format, ignoring that they both start at overlapping offsets, and keeping the head.



Just parasitize `FileA` with `BodyB` and keep `HeadB` as the overlap:



The minimum length overlap is basically the number of bytes where you can declare a file then have some unparsed space, either naturally or by declaring a comment. This value can change drastically between file formats, as shown on page 15.

For example, it's 1 for PostScript because you can declare a line comment with `%`, so provided there's no encoded newline *after* decryption, this will be a valid ignored space. The PE file format's minimum overlap is two bytes — `MZ` — because you can abuse the DOS header, limiting you to a 58-byte parasite.

A JPEG header with comment declaration is `FF D8 FF FE XX YY`, which is six bytes, with `XXYY` being the length of the comment in big endian. However, if you need a `0x3489`-long comment, a `0x35??`-long comment will do the trick, so you don't have to bruteforce the `YY` byte. If you feel luck, you might

Minimal start offset	Variable offset												Unsupported parasite																					
	1	2	4	8	9	16	20	23	28	34	40	64	94	132	12	28	16																	
	P	P	J	F	M	T	F	W	G	P	R	I	R	B	C	I	P	C	J	P	E	A	P	I	I	J	W	B	O	B	E	G	L	N
	S	E	P	l	P	I	L	A	Z	N	I	D	T	M	P	L	S	A	P	C	L	R	C	C	C	a	A	P	G	Z	B	I	N	E
			G	a	4	F	V	D		G	F	3	F	P	I	D	D	B	2	A	F		A	O	C	v	S	G	G	2	M	F	K	S
			c	F					F	v					O	A		P		P		a	M					L						
1* PS	.	M	A	?	?	?	?	?	A	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?	?
2@ PE	M	.	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	!	!	!	!	!	M	M	M	!	!	!	!	!	
4+ JPG	A	A	.	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	
. .	[the table could go on but would take too long to bruteforce]																																	
	X:	automated	?:	likely possible																														
	M:	manual	!:	unknown																														

- * Hack that relies on line comments with GhostScript. Requires the parasite not to contain any new line, after encryption.
- @ Hack relying on overwriting the DOS Header, therefore restricting the parasite space to offsets 2-60.
- + Signature, comment declaration and length are two bytes apiece. To specify them all is six bytes, but if we round up the big-endian length and leave its low byte uncontrolled, we only need five. And if we leave the length entirely uncontrolled, we only need to fix four bytes.

Figure 3: Minimal start offsets of file formats, and exploitation via near polyglots

The Inforex 180 Magnetic Line Printer.

Only the paper is ordinary.

Our patented dry-ink transfer gives you a high quality printout on an ordinary, inexpensive 8½" roll of paper. But everything else about the 180 is definitely *not* ordinary.

QUIET.
At least 10db below electric typewriters. Ideally suited for the office environment.

SMALL.
Desk top size, weighs just 33 pounds.

INTERFACE.
Serial RS-232 or TTL. Parallel TTL.

CHARACTER SET.
Full 96 characters ASCII, upper and lower case. Expanded character sets available as options.



FAST.
180 lines per minute. A 1920 character screen in 8 seconds.

GRAPHICS OPTION.
Permits intermix of text with bar charts, curves, etc.

APPLICATIONS.
CRT hard copy, minicomputers, remote printing operations, OEM or end-user.

INFOREX
For more information, write Inforex, Incorporated, Dept. 588 CD-7, 21 North Ave., Burlington, MA 01803

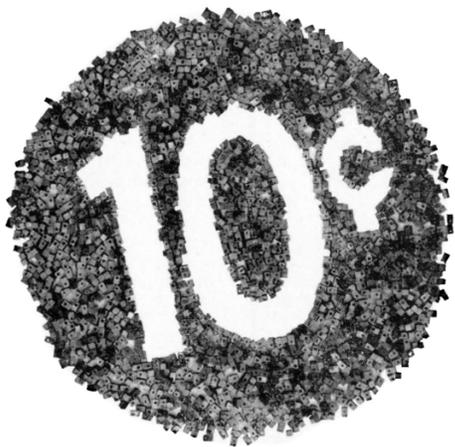
also gamble on the length and not bother to brute-force XX either.

Mitra can generate such files with the `--overlap` parameter. It keeps the overlap's content in the filename as well as the offsets where the content changes formats, to be re-used later by AngeCryption or TimeCryption scripts.

Ambiguity Files with different interpretations depending on parsers we now call *ambiguous* and previously called *schizophrenic*, werewolves or shapeshifters. There are plenty of sources of this ambiguity.

When a value such as a pointer never changes across standard files, it's tempting for a parser to simply ignore it. Putting some contents under unusual conditions while putting other contents under the typical conditions might reveal a difference between the two parsers.

Sometimes a value is represented twice. For example, a buffer with a declared length might also end with a null terminator. What if that terminator happens earlier than the declared length? Which length value is the real one? Or if you declare the same value twice, and there isn't an error, does the first or the second declaration take priority?



BIAX® MEMORIES NOW 10¢ A BIT

The new BIAx NANOLOK electrically alterable memory gives 3-Plus megacycle TRUE NDRO operation at the low cost of 10 cents per bit in 200,000 bit systems in quantity. NANOLOK is designed for commercial and industrial data use — but can be adapted readily to shipboard and mobile MIL-Spec environments. BIAx used to be expensive, but not any more — not with NANOLOK. Do your bit and write today for all the details on NANOLOK in Data File B-129.



RAYTHEON COMPUTER, 2700 South Fairview Street, Santa Ana, California 92704

If you corrupt a format on purpose and the parser tries to rebuild the file, how does it do it first? What if you put a valid file structure in a comment? Such recovery algorithms are typically not officially specified, so each developer might do it differently.

Some formats are extending older formats. Both the old and the new formats are present in the file. These formats are naturally ambiguous at a format level, and we might call them ambiguous polyglots.

A widespread example is the Portable Executable, defined as an extension of the DOS format. Preciously few PE files—such as `regedit95.exe`—have a meaningful DOS payload. Valid PE files are expected to just have the same DOS stub with no unique code.

Robert Xiao proved us wrong by crafting a universal Doom binary, which works from DOS 6 to Windows 10, both as valid DOS and PE payloads in the same file.⁹ This is something like a father and son having the same name, with no distinctive suffix whatsoever.

Ciphertexts can be ambiguous too, even despite authenticated encryption! Check the TimeCryption article page 30 in this issue for abuses of GCM, GCM-SIV and OCB3.

Collisions PoC||GTFO 19:05 covered a lot of details for exploiting hash collisions. However `tar.gz` and `DOCX` (ZIPped XML)—which were initially thought to be unexploitable—are properly dealt with and explained in the Inside Out article on page 19.

Conclusion

With basic knowledge of file format identification and abuse, Mitra can try different strategies and generates many forms of file abuse: payload embedding, mock files, polyglots and pseudo-polyglots.

Pseudo-polyglots are the unified. form of file formats abuses to be combined with cryptographic operations. They include both AngeCryption, covering ECB, CBC, CFB, OFB modes, and TimeCryption, covering CTR, OFB, GCM, OCB3, GCM-SIV modes.

Extensions of Mitra might cover ambiguous files with standard strategies, hash collisions and hash collisions over different formats.

⁹[git clone https://github.com/nneonneo/universal-doom](https://github.com/nneonneo/universal-doom)