



# The Treachery of Files and Two New Tools that Tame It

Evan Sultanik

# The Treachery of Files

A Tragedy in Two Acts

---

## Act I

Files Have No Intrinsic Meaning

**Goal:** Convince you that these funky files aren't just nifty parlor tricks

## Act II

PolyFile and PolyTracker!

**Goal:** Introduce two new tools to help you reverse engineer files and the parsers that process them

# whoami

Résumé: Evan A. Sultanik, Ph.D.

November 29<sup>th</sup>, 2017

evan@sultanik.com  
[https://sultanik.com/](https://sultanik.com) <https://github.com/ESultanik> <https://twitter.com/ESultanik> <https://keybase.io/ESultanik>

2017–Present	Trail of Bits	Security Researcher
2012–2017	Digital Operatives	Chief Scientist
2010–2012	The Johns Hopkins University APL	Senior Research Scientist
2006–2015	Drexel University Department of Comp. Sci.	Adjunct Faculty
2001–2010	Drexel University Department of Comp. Sci.	Research Fellow
1998–2003	Feith Systems and Software	Software Engineer

The following is a chronological account of some of the things I've done and software I've created recently. I am legally constrained about talking about portions of my work (*e.g.*, work for certain government customers, code audits, and expert witness for IP litigation), so much of it is omitted here. In lieu of a traditional laundry list of technologies and tools with which I'm familiar, I hope the following will serve to demonstrate that I am a generalist with strong computer science fundamentals and the ability to quickly learn and master new technologies. I specialize in computer security, AI/ML/NLP, combinatorial optimization, and distributed systems. You might also note that I've unabashedly written this informally. Do you trust people who speak about themselves in the third person? I don't. For a more traditional résumé, my formal academic CV is available on my website.

**August, 2017** Edited and co-published *PoC||GTFO*, a nearly 800 page print collection gathering articles published in the first ten issues of the International Journal of Proof-of-Concept or GTFO. PoC||GTFO—on which I have been an editor since August of 2015 and a contributor since June of 2014—follows in the tradition of *Phrack* and *Uninformed* by publishing on the subjects of offensive security research, reverse engineering, and file format internals. Until now, the journal has only been available online or printed and distributed for free at conferences worldwide. ISBN-13: 978-1-59327-880-9 <https://www.nostarch.com/gtfo>

**October, 2015–August, 2017** Member of the “control team” of expert *human* vulnerability researchers for the DARPA VET program, the purpose of which was to act as a baseline against which to compare the *automated* vulnerability analysis tools of the other performers. This consisted of using state-of-the-art publicly available, commercial-, and government-off-the-shelf tools (*e.g.*, IDA) to determine whether a given set of black-box ARM binaries had been modified to include malice. Our team outperformed the automated equivalents across four judged engagements. Primarily IDA, HexRays, and our own static analysis tools that I helped develop (see below).

**June, 2017** Published and presented an analysis of zoning and residential density changes in Philadelphia over the past five years. This required devising and implementing a solution to the difficult computational geometry problem of intersecting the tens of thousands of latitude/longitude polygons describing Philadelphia's zoning districts both pre and post the comprehensive zoning overhaul that occurred in 2012. Also used Machine Learning to create a model to predict the real estate tax revenue implications of zoning changes. Implemented in Python, Java, and KML. <https://www.sultanik.com/blog/ZoningDensity>

**June, 2017** Created a PDF that is also a valid Git repository containing its own L<sup>A</sup>T<sub>E</sub>X source code and a copy of itself<sup>1</sup>. Published in PoC||GTFO Issue 0x15. Implemented in Bash, L<sup>A</sup>T<sub>E</sub>X, Python, and C. <https://github.com/ESultanik/PDFGitPolyglot>

**June, 2017** Created a toy web-based terminal emulator and shell, backed by a filesystem served over HTTP (similar to WebDAV), in which scripts are dynamically loaded from the server via JavaScript, and replete with libraries like ncurses. Implemented purely in JavaScript using JQuery. <https://www.sultanik.com/#term>

**March, 2017** Created a PDF that is also a valid Nintendo Entertainment System ROM image<sup>2</sup> that, when played/emulated, renders the title page of the PDF and displays the MD5 checksum of itself<sup>3</sup>. This MD5 quine is achieved by solving for 128 MD5 collisions using an implementation of an algorithm created by Marc Stevens. This polyglot was in collaboration with Evan Teran and Ange Albertini. Implemented in a mixture of C, 6502 assembly, and Python. <https://www.sultanik.com/pocorgtfo/#0x14>

**August, 2013–November, 2016** Created a semi-automated static analysis tool for human-in-the-loop reverse engineering and vulnerability research along with four colleagues. Led the effort in researching and implementing components for static taint analysis, and devised novel algorithms for SMT constraint optimization and detecting algorithmic complexity defects. Implemented in C++11.

**October, 2016** Discovered a vulnerability in the way most US defense contractors choose passwords for the spinlocks they are required to use to secure safes and rooms containing classified information. I wrote a short closet drama about it. Implemented in Bash, Python, and English Prose. <https://archive.org/stream/pocorgtfo13#page/n42/mode/1up>

**October, 2016** Created a PDF that is also a valid PostScript file. If you send the raw PDF to a printer, it will print out differently than how it was rendered on screen in the PDF viewer. For good measure, the PostScript also reads your /etc/passwd file. Implemented in Python and PostScript. <https://www.sultanik.com/pocorgtfo/#0x13>

**August, 2016** Along with three collaborators, created an AI bot capable of playing Pokémon Go. Unlike other bots of the time which exploited bugs in Niantic's implementation to achieve superhuman feats like teleportation to increase their speed,

<sup>1</sup>unzip ESultanikResume.pdf ESultanikResume/PDFGitPolyglot.pdf

<sup>2</sup>This PDF is also a valid Nintendo Entertainment System ROM<sup>3</sup>. Try emulating it!

<sup>3</sup>By the way, the NES ROM version of this PDF also prints out its MD5.

<sup>4</sup>By the way, the MD5 hash of this PDF is AFFFAB1EDA8E9B6D8A80E940F20CB3B3B<sup>5</sup>.

<sup>5</sup>By the way, the first seven characters of that MD5 hash are not a coincidence, either.

# `whoami`

---

<sup>1</sup>`unzip ESultanikResume.pdf ESultanikResume/PDFGitPolyglot.pdf`

<sup>2</sup>This PDF is also a valid Nintendo Entertainment System ROM<sup>3</sup>. Try emulating it!

<sup>3</sup>By the way, the NES ROM version of this PDF also prints out its MD5.

<sup>4</sup>By the way, the MD5 hash of this PDF is `AFFAB1EDA8E9B6D8A80E940F20CB3B3B`.<sup>5</sup>

<sup>5</sup>By the way, the first seven characters of that MD5 hash are not a coincidence, either.

# `whoami`



---

<sup>1</sup>`unzip ESultanikResume.pdf ESultanikResume/PDFGitPolyglot.pdf`

<sup>2</sup>This PDF is also a valid Nintendo Entertainment System ROM<sup>3</sup>. Try emulating it!

<sup>3</sup>By the way, the NES ROM version of this PDF also prints out its MD5.

<sup>4</sup>By the way, the MD5 hash of this PDF is `AFFAB1EDA8E9B6D8A80E940F20CB3B3B`.<sup>5</sup>

<sup>5</sup>By the way, the first seven characters of that MD5 hash are not a coincidence, either.

```
$ md5sum ESultanikResume.pdf  
affab1eda8e9b6d8a80e940f20cb3b3b ESultanikResume.pdf
```

```
$ md5sum ESultanikResume.pdf  
affab1eda8e9b6d8a80e940f20cb3b3b ESultanikResume.pdf
```

```
$ unzip -l ESultanikResume.pdf
```

```
Archive: ESultanikResume.pdf
```

Length	Date	Time	Name
-----	-----	-----	-----
0	09-17-2019	14:47	ESultanikResume/
638270	06-07-2019	10:35	ESultanikResume/PDFGitPolyglot.pdf
-----	-----	-----	-----
638270			2 files

```
$ md5sum ESultanikResume.pdf
affab1eda8e9b6d8a80e940f20cb3b3b  ESultanikResume.pdf
$ unzip -l ESultanikResume.pdf
Archive:  ESultanikResume.pdf
      Length      Date  Time    Name
-----  -----
          0  09-17-2019 14:47  ESultanikResume/
  638270  06-07-2019 10:35  ESultanikResume/PDFGitPolyglot.pdf
-----  -----
      638270
$ file ESultanikResume.pdf
ESultanikResume.pdf: NES ROM image (iNES): 8x16k PRG, 2x8k
CHR [V-mirror] [SRAM] [Trainer]
```

```
$ md5sum ESulta
affab1eda8e9b
$ unzip -l ES
Archive:  ESu
Length
-----
0 09
638270 06
-----
638270
$ file ESulta
ESultanikResu
CHR [V-mirror]
$ nestopia ES
```

```
df
tPolyglot.pdf
G, 2x8k
```

# Act I

Files Have No Intrinsic Meaning

PoC||GTFO

<https://sultanik.com/pocorgtfo/>

# PoCorGTF0

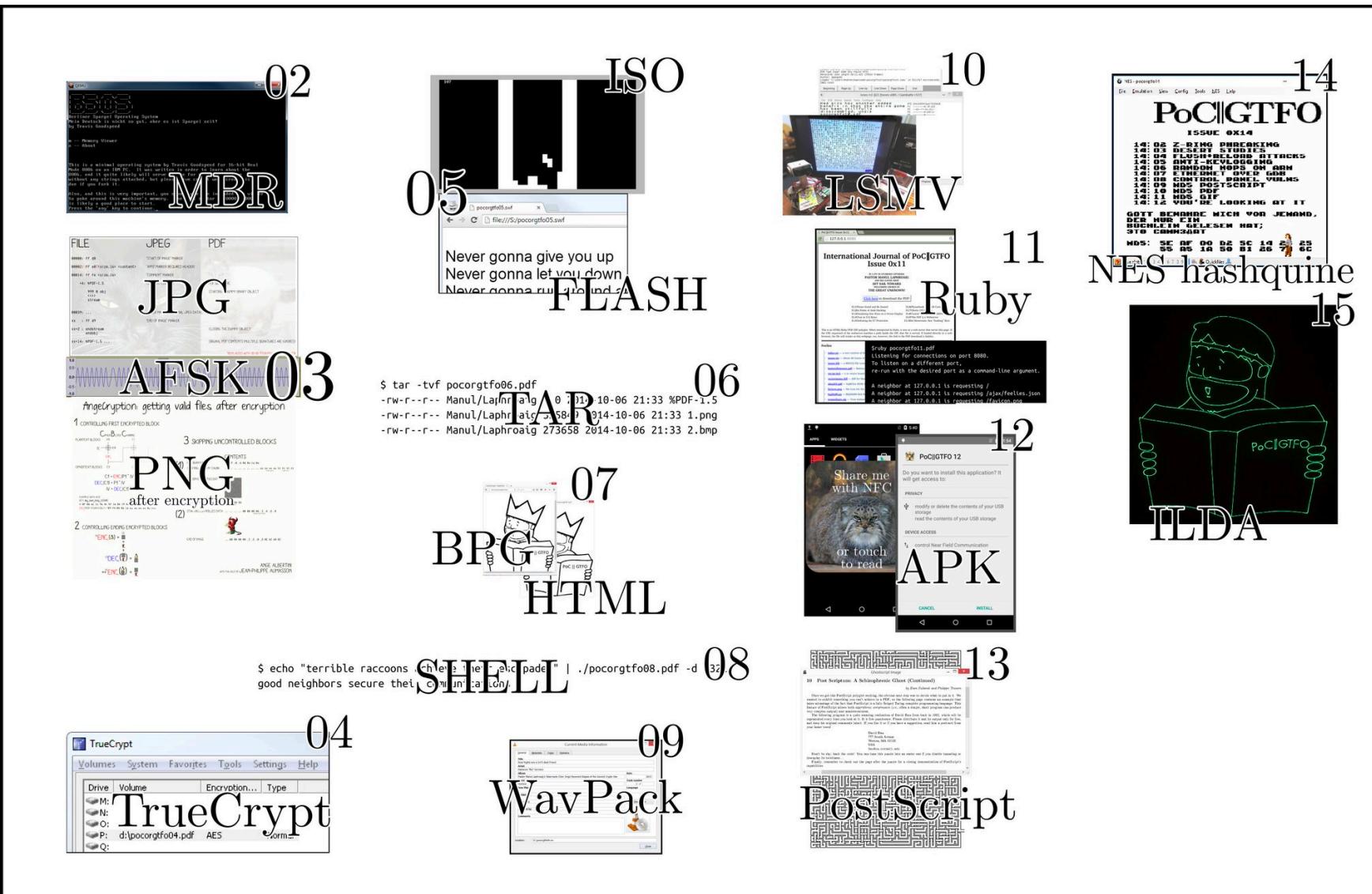
# Proof of Concept

## (Pictures of Cats)

*“It looks great on a shelf, and if you read PoC//GTFO on public transportation, people stay away from you.”*

—Hackaday Review

Roughly quarterly journal, in the tradition of Phrack and Uninformed  
Offensive security research and stunt hacking  
First released on paper at a conference, later released digitally  
Each digital release is a Polyglot



<https://sultanik.com/pocorgtfo/>

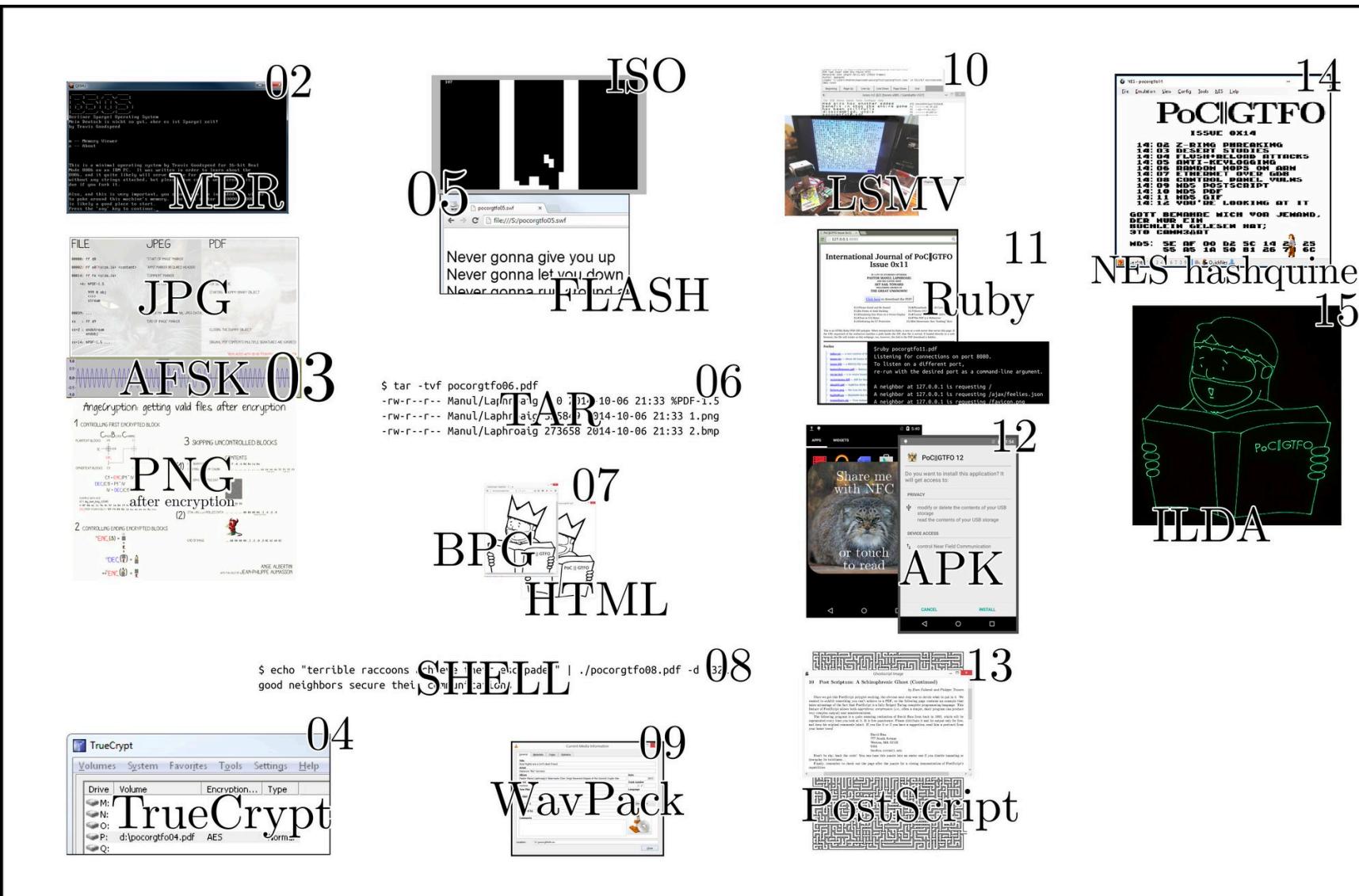
# PoCorGTFO

Proof of Concept  
(Pictures of Cats)

*“It looks great on a shelf, and if you read PoC//GTFO on public transportation, people stay away from you.”*

—Hackaday Review

Roughly quarterly journal, in the tradition of Phrack and Uninformed  
Offensive security research and stunt hacking  
First released on paper at a conference, later released digitally  
Each digital release is a Polyglot



**Neil Madden** @neilmaddog · Jul 19

I wonder what happened to PoC||GTFO issues 0x0A–0x0F...



**Evan Sultanik**

@ESultanik

Replying to @neilmaddog

We number in BCD in honor of the HP48 calculator's floating point implementation, which matches decimal rounding errors.

3:01 PM - 19 Jul 2017

<https://sultanik.com/pocorgtfo/>

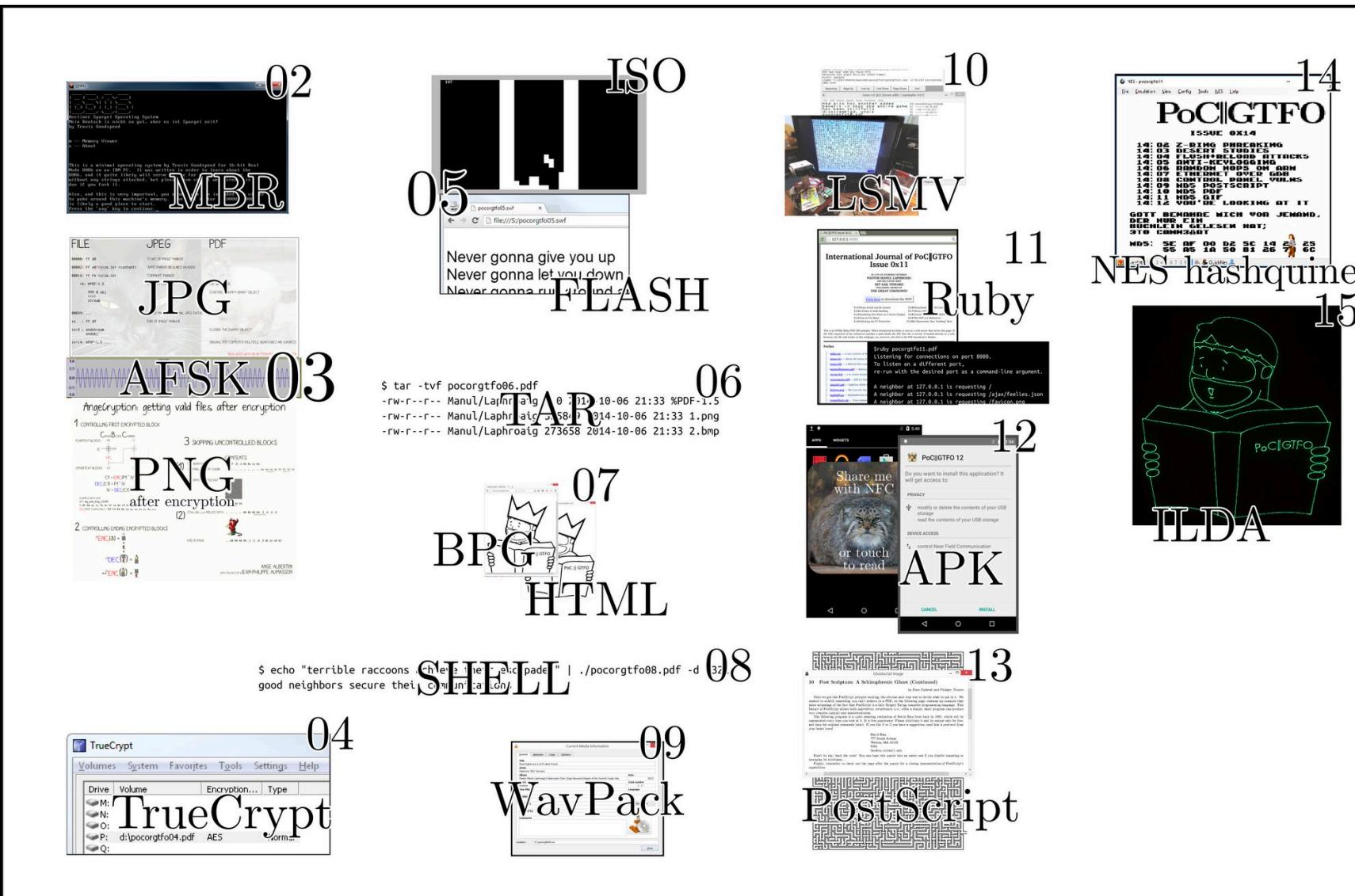
# PoCorGTFO

Proof of Concept  
(Pictures of Cats)

*“It looks great on a shelf, and if you read PoC//GTFO on public transportation, people stay away from you.”*

—Hackaday Review

Roughly quarterly journal, in the tradition of Phrack and Uninformed  
Offensive security research and stunt hacking  
First released on paper at a conference, later released digitally  
Each digital release is a Polyglot



**Neil Madden** @neilmaddog · Jul 19

I wonder what happened to PoC||GTFO issues 0x0A–0x0F...



**Travis Goodspeed**

@travisgoodspeed

Replying to @\_gbg\_ @h2hconference @hacktivityconf

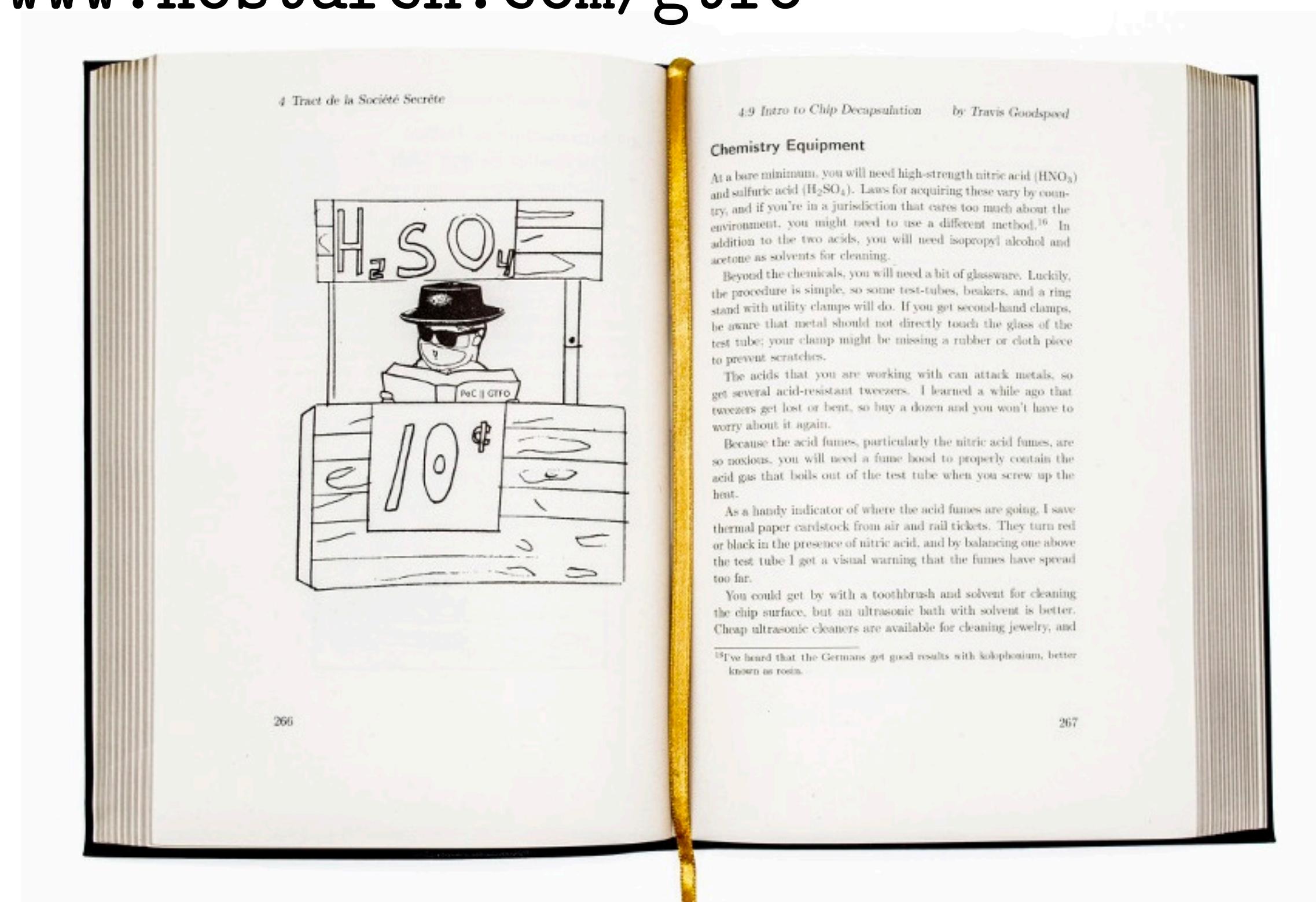
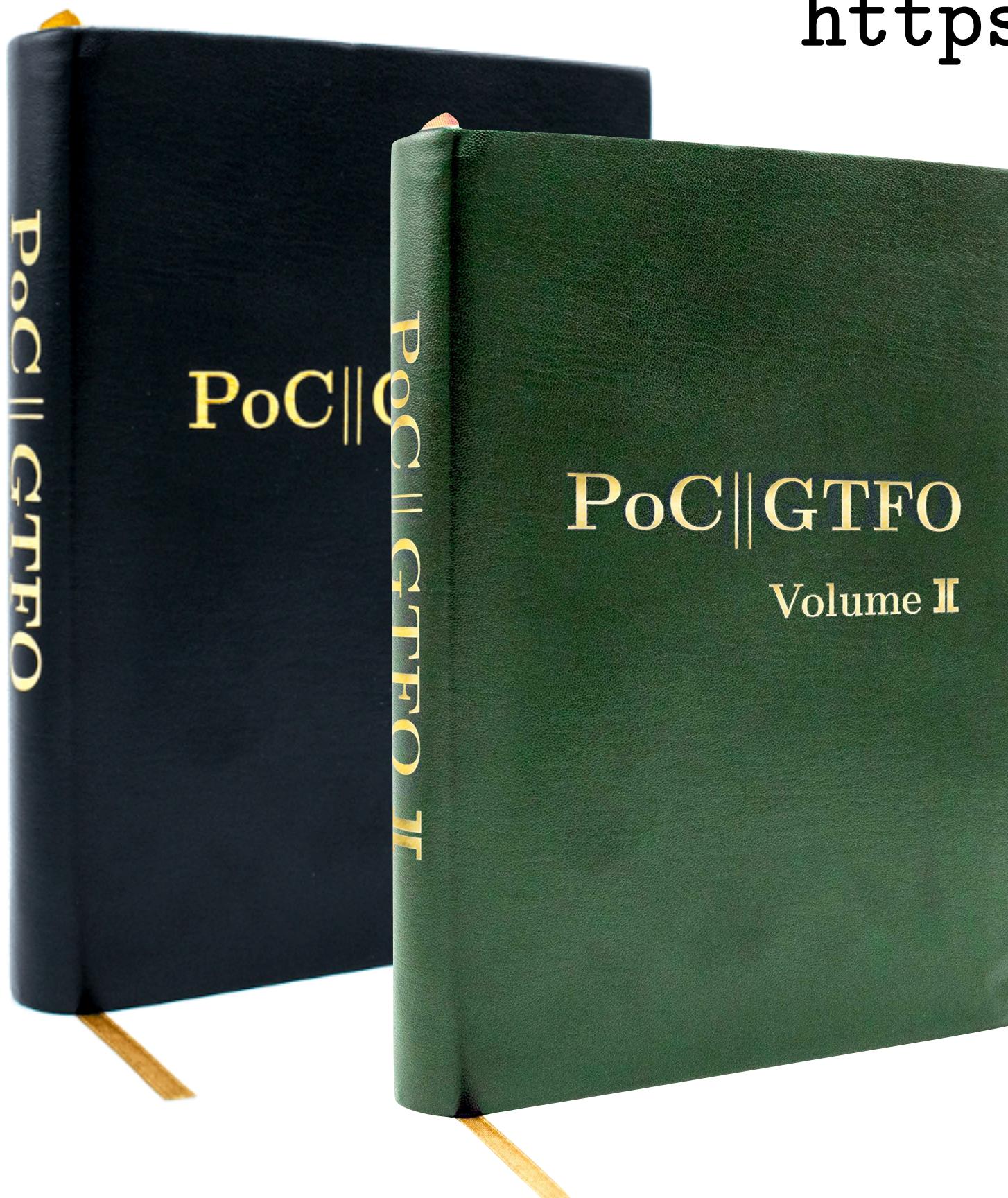
We number them in binary coded decimal, as  
a tribute to the floating point unit of the HP  
SATURN architecture.

1:16 PM - 12 Oct 2017

<https://sultanik.com/pocorgtfo/>

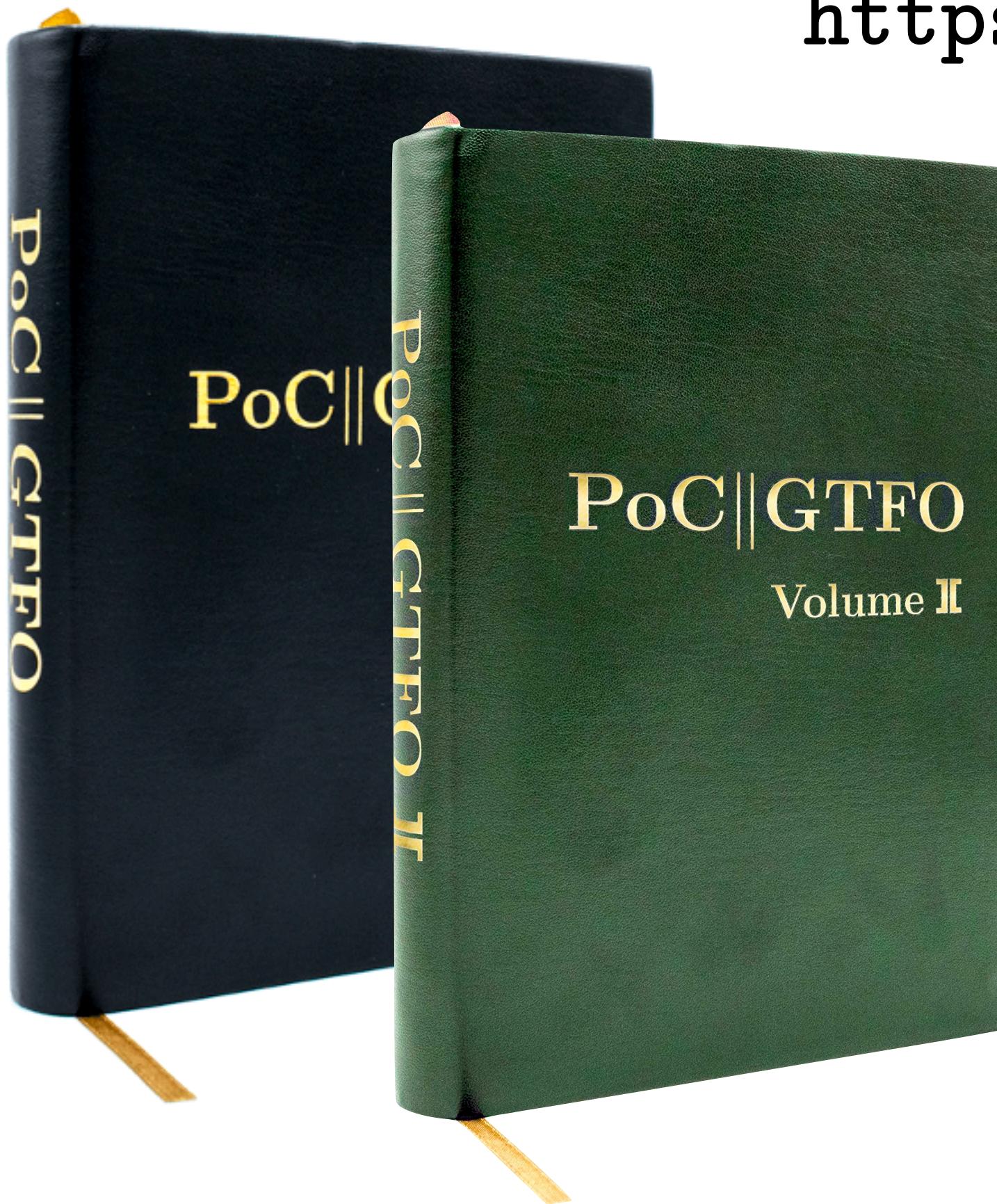
# The Book of PoC||GTFO

<https://www.nostarch.com/gtfo>



# The Book of PoC||GTFO

<https://www.nostarch.com/gtfo>

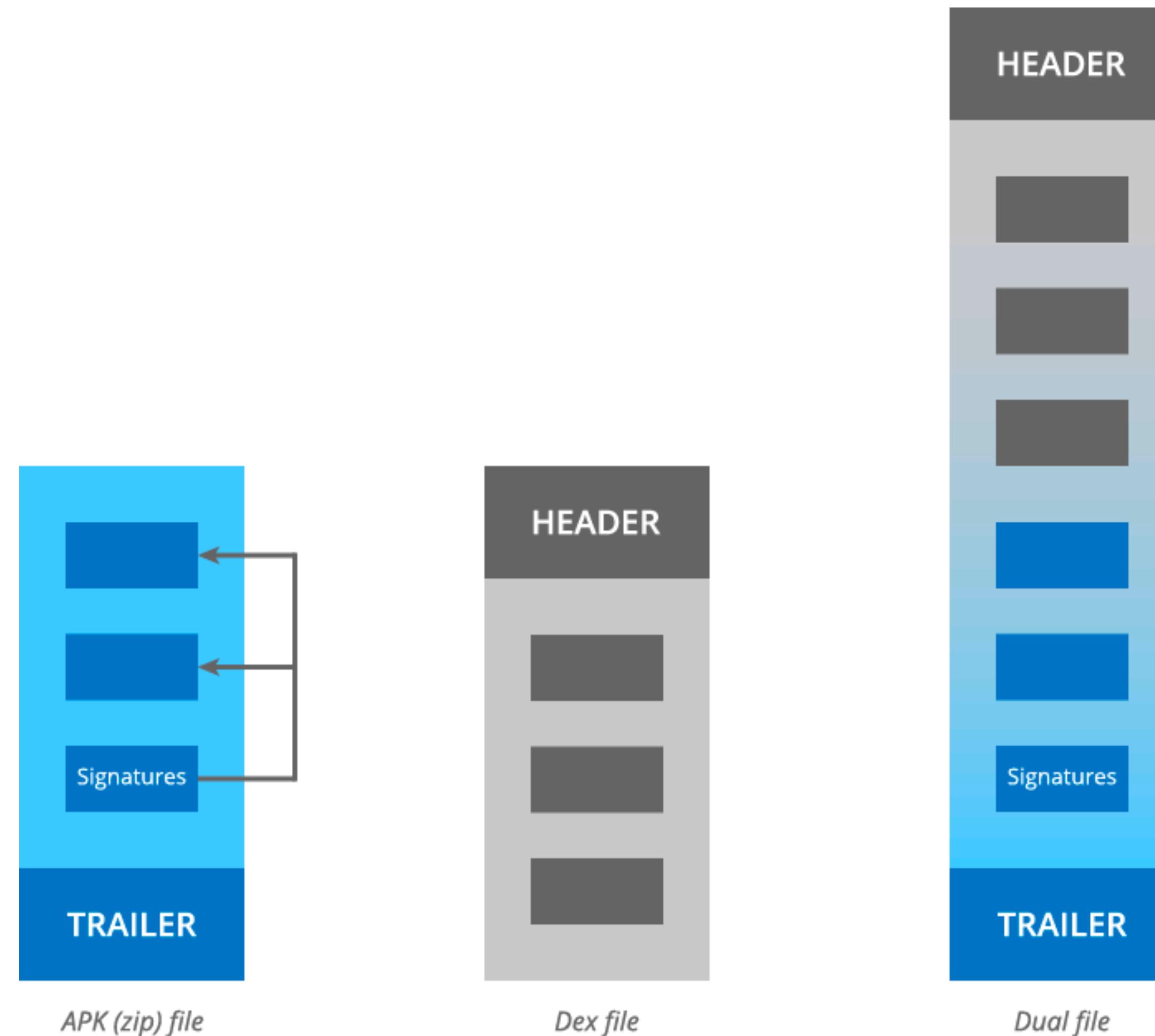


*“...a file has no intrinsic meaning. The meaning of a file—its type, its validity, its contents—can be different for each parser or interpreter”*

—PoC||GTFO 7:6 by Ange Albertini

# Example: Android

## APK (zip)/Dex Polyglot



# Example: Android APK (zip)/Dex Polyglot



APK (zip) file



# Example: Android APK (zip)/Dex Polyglot



July, 2017



# Example: Android APK (zip)/Dex Polyglot



July, 2017



# What's Wrong Here?

```
$ tar xvf totally_not_malware.tar.gz
```

# What's Wrong Here?

```
$ tar xvf totally_not_malware.tar.gz
```



Note: We didn't provide the z option!

Modern versions of tar automagically detect that the archive is compressed based on magic bytes!  
(The actual file extension is ignored.)



# What's Wrong Here?

```
$ tar xvf totally_not_malware.tar.gz
```



Note: We didn't provide the z option!

Modern versions of tar automagically detect that the archive is compressed based on magic bytes!  
(The actual file extension is ignored.)



What if we created a file that is *both* a valid .tar and a valid .tar.gz?



`totally_not_malware.tar`





# Why are PDFs Particularly Polyglottable?

- Because “Adobe,” that’s why!
- It’s been around for a long time
- Parsers built to be resilient to all sorts of errors and incompatibilities
- Can insert arbitrary length binary blobs almost anywhere in the file
- Almost all parsers ignore everything before the header

```
9999 0 obj
<<
/Length # bytes in the blob
>>
stream
lol, put whatever you want here!
endstream
endobj
```

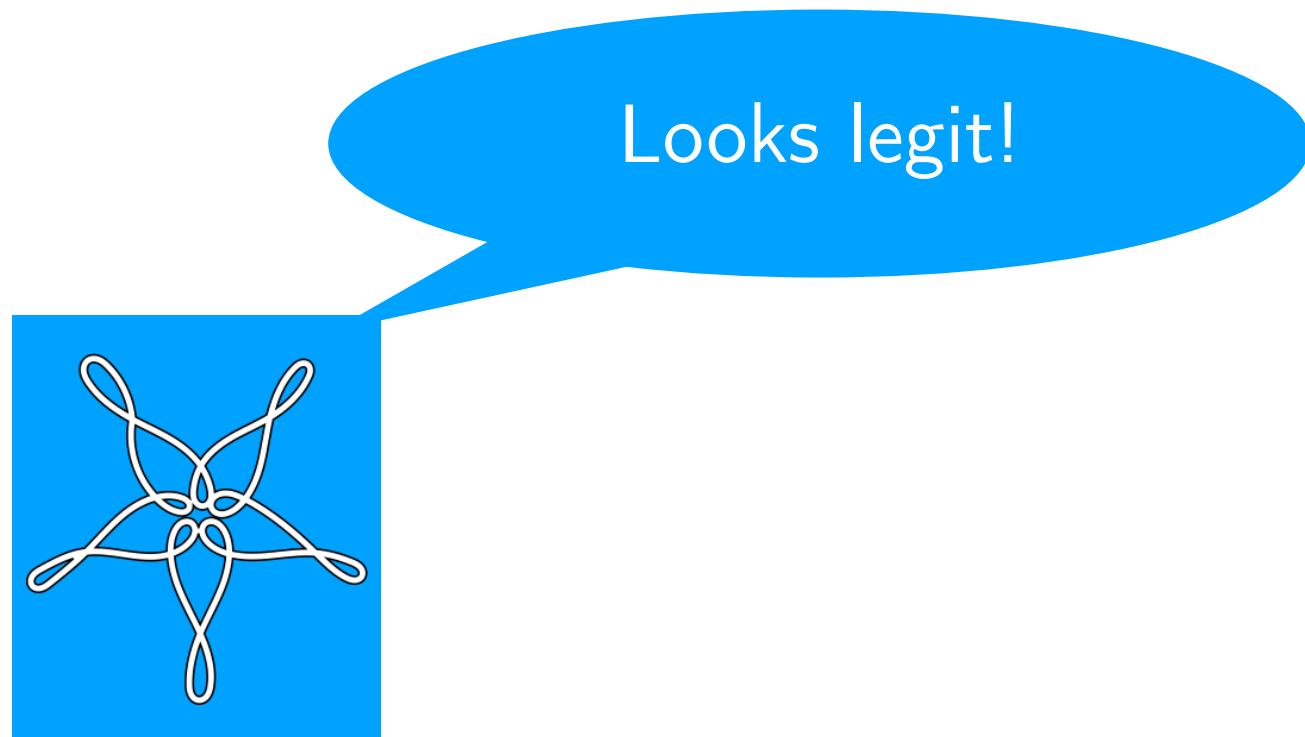
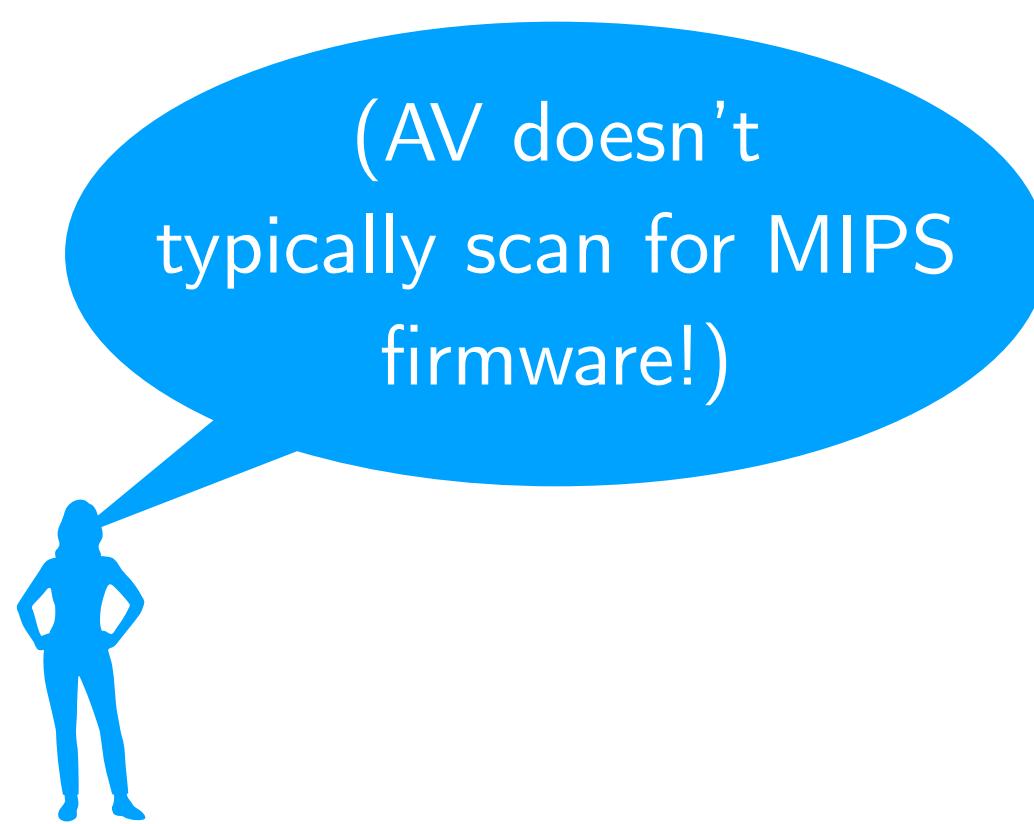
```
lol, put whatever you want here!
%PDF-1.5
%<D0><D4><C5><D8>
:
```



Hey,  
I've been trying to get my résumé to  
so-and-so in HR, but we've had problems  
with E-mail. Can you please print out the  
attached copy and give it to them?

Thanks! —Alice Hackerman

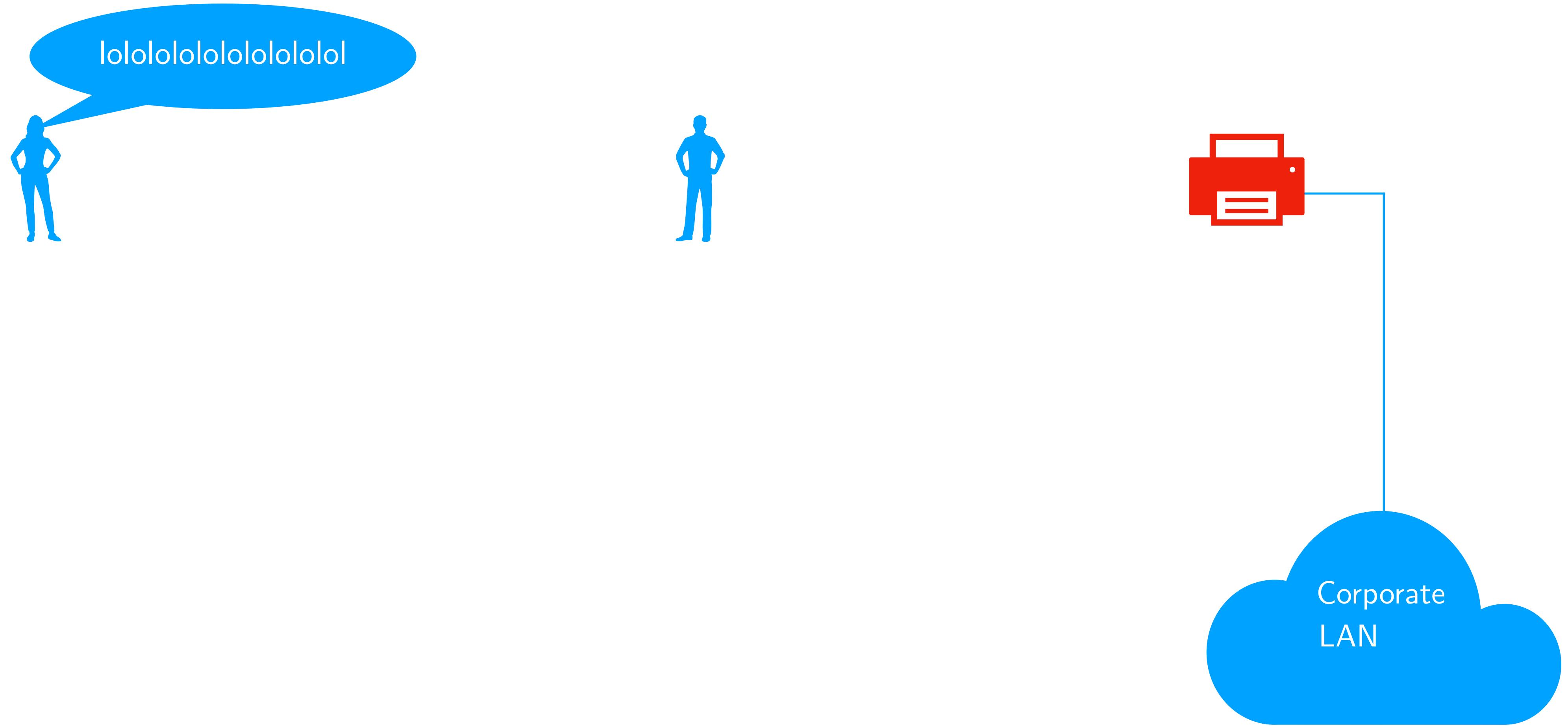




Hey,  
I've been trying to get my résumé to  
so-and-so in HR, but we've had problems  
with E-mail. Can you please print out the  
attached copy and give it to them?

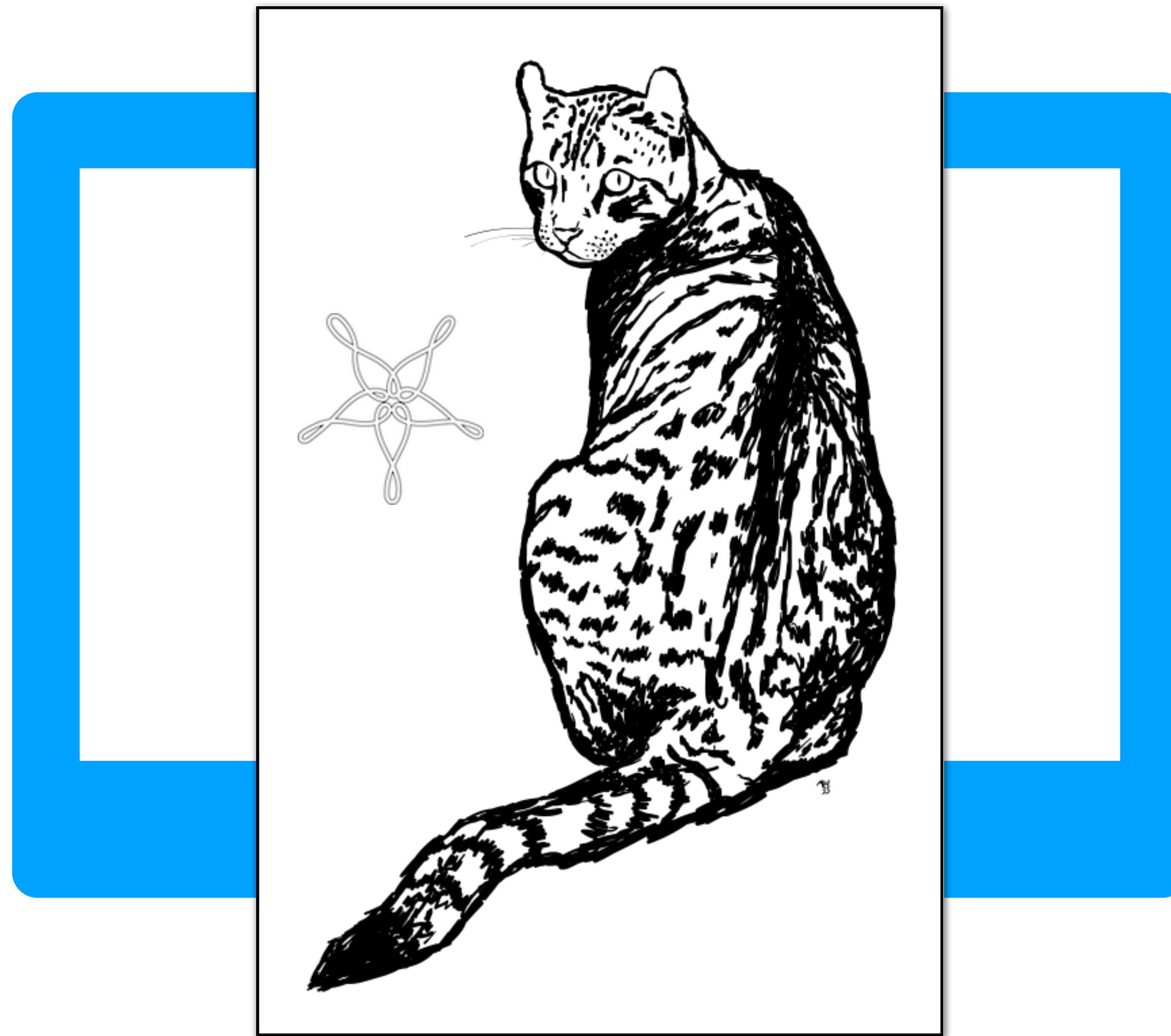


Thanks! —Alice Hackerman

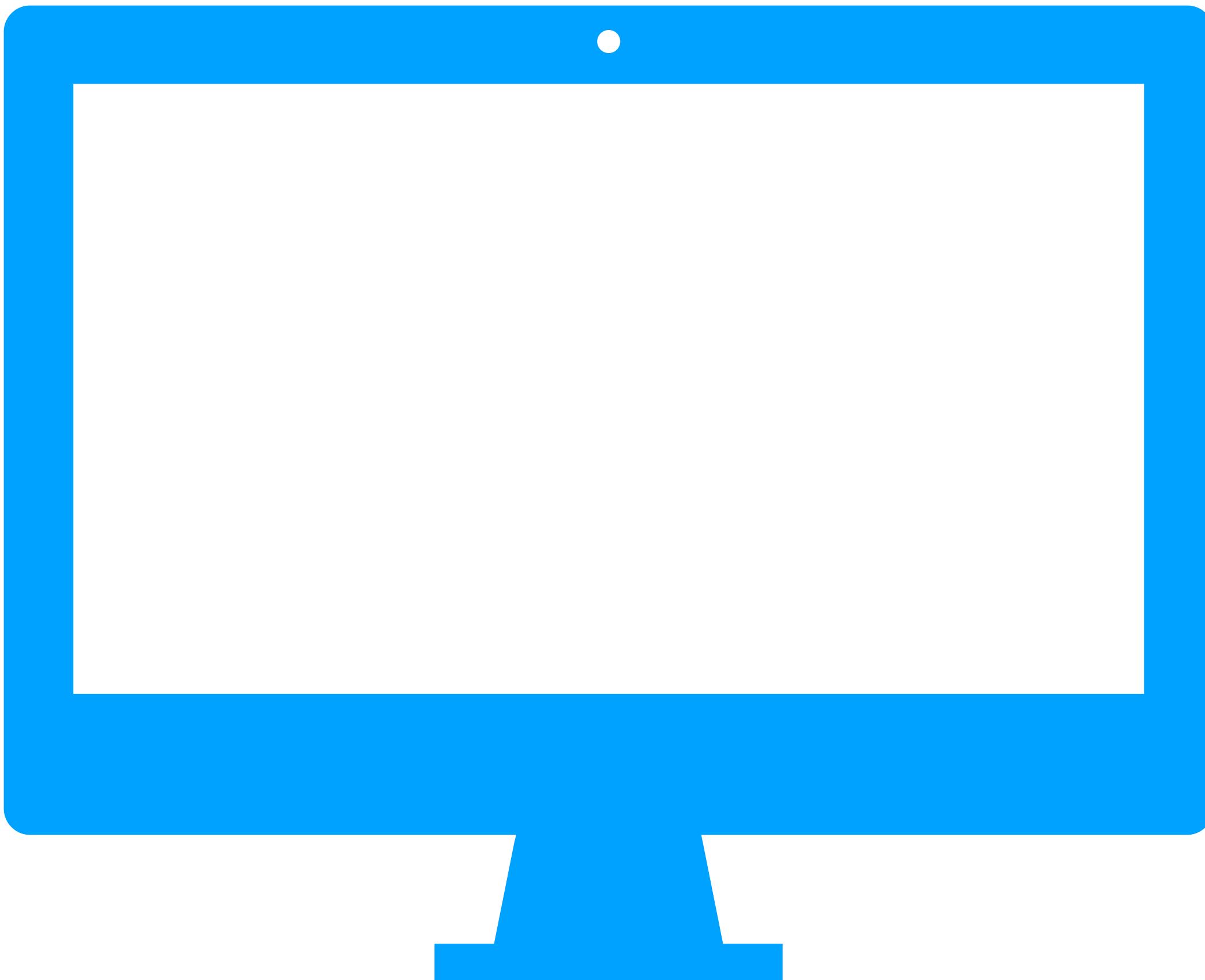


PostScript/PJL (Cui & Stolfo, 2011)

# Producing a Positively Provocative PDF/PostScript Polyglot

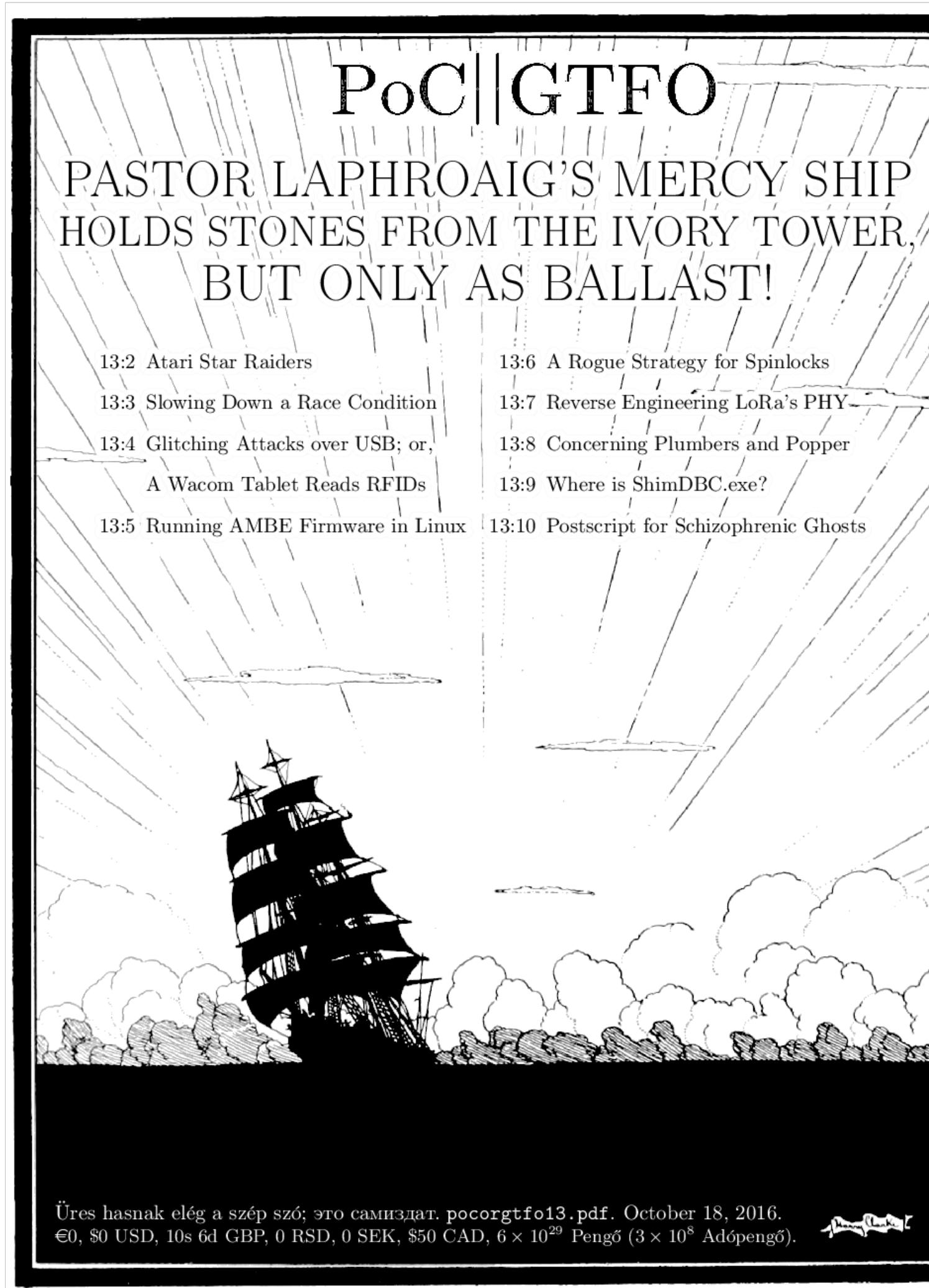


# Producing a Positively Provocative PDF/PostScript Polyglot



**Don't print  
PostScript  
created  
by Evan!**

# PoC||GTFO 0x13



PDF

## 10 Post Scriptum: A Schizophrenic Ghost (Continued)

by Evan Sultanik and Philippe Teuwen

Once we got this PostScript polyglot working, the obvious next step was to decide what to put in it. We wanted to exhibit something you can't achieve in a PDF, so the following page contains an example that takes advantage of the fact that PostScript is a fully fledged Turing complete programming language. This feature of PostScript allows both *algorithmic compression* (i.e., often a simple, short program can produce very complex output) and nondeterminism.

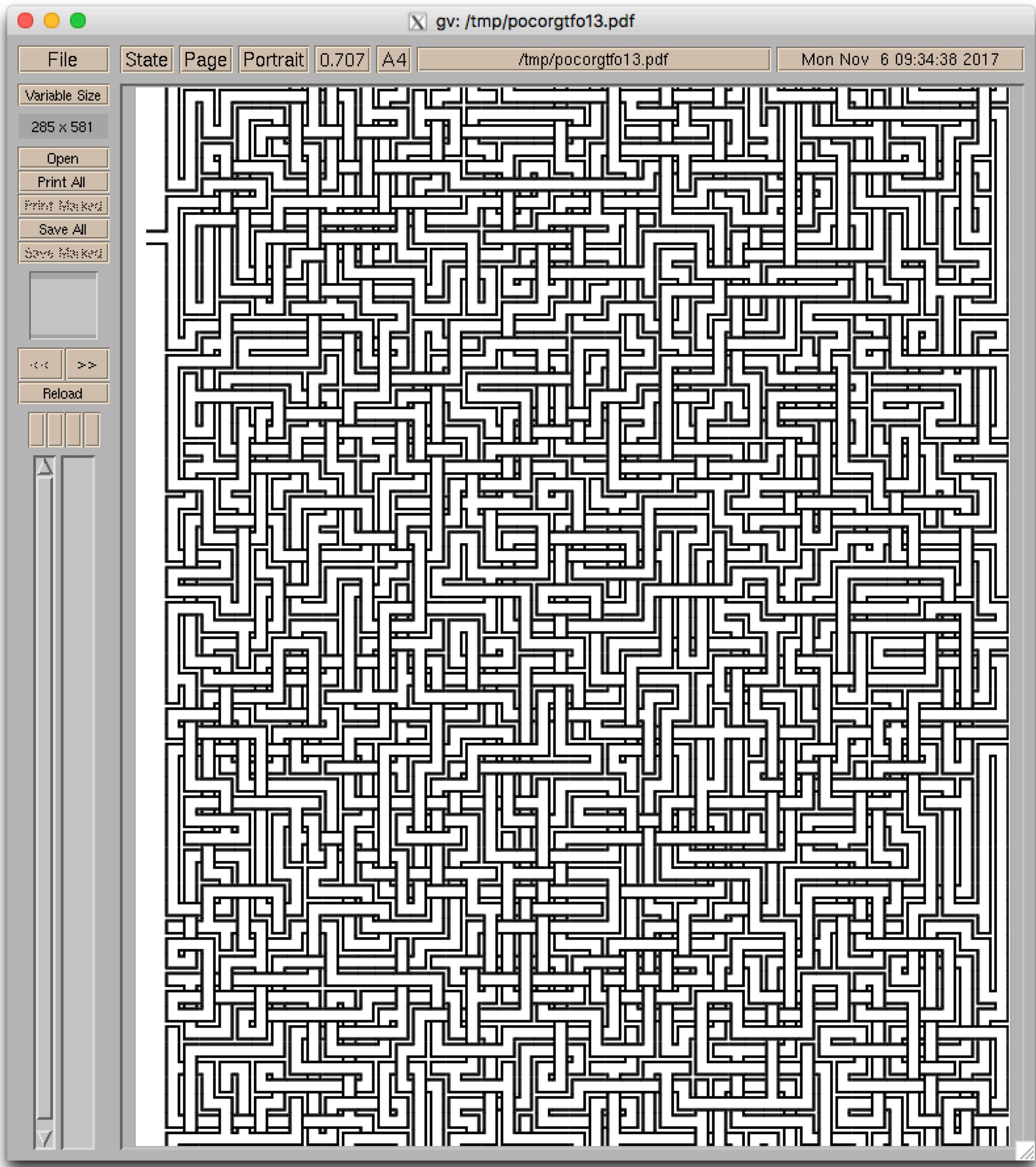
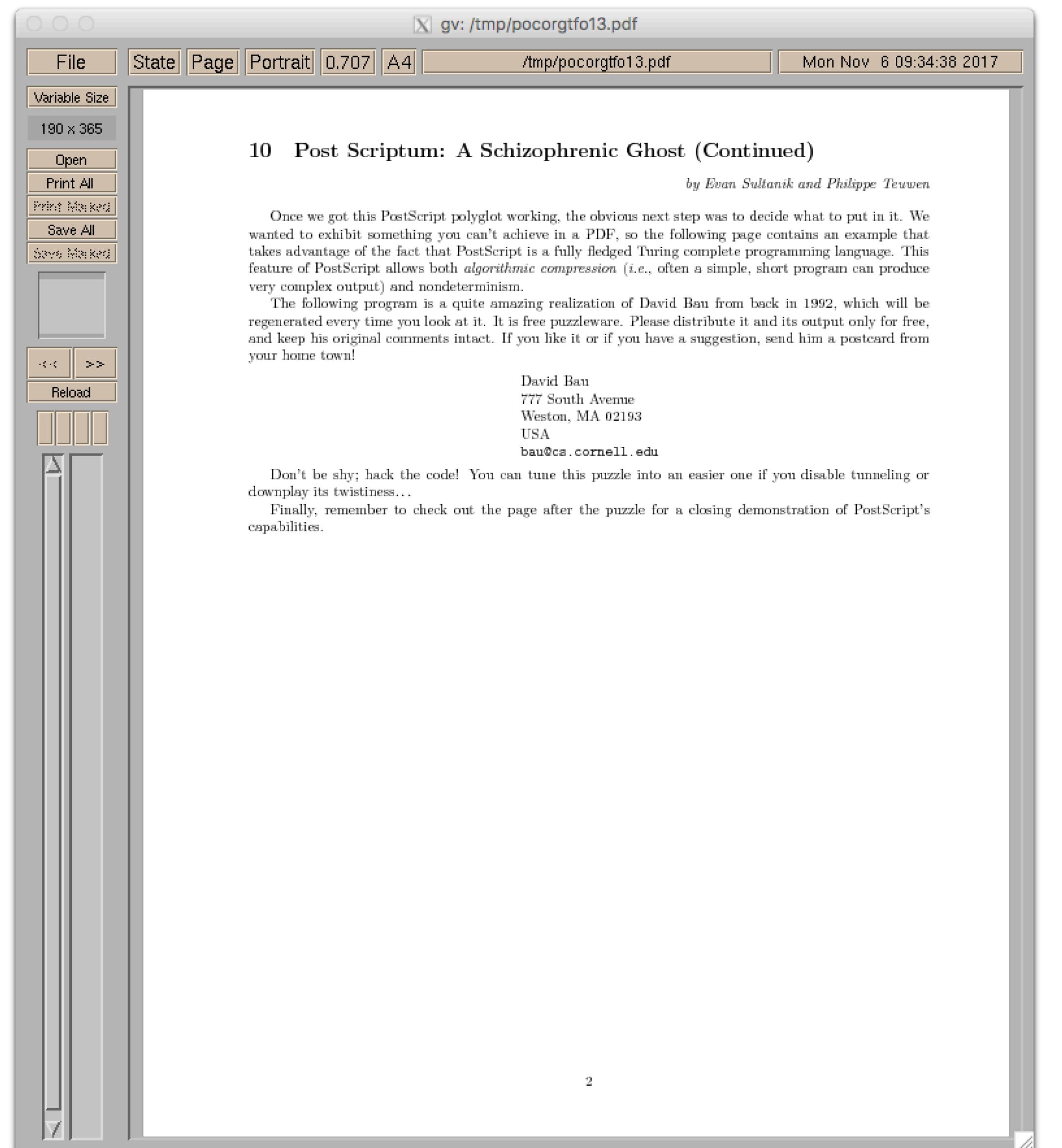
The following program is a quite amazing realization of David Bau from back in 1992, which will be regenerated every time you look at it. It is free puzzleware. Please distribute it and its output only for free, and keep his original comments intact. If you like it or if you have a suggestion, send him a postcard from your home town!

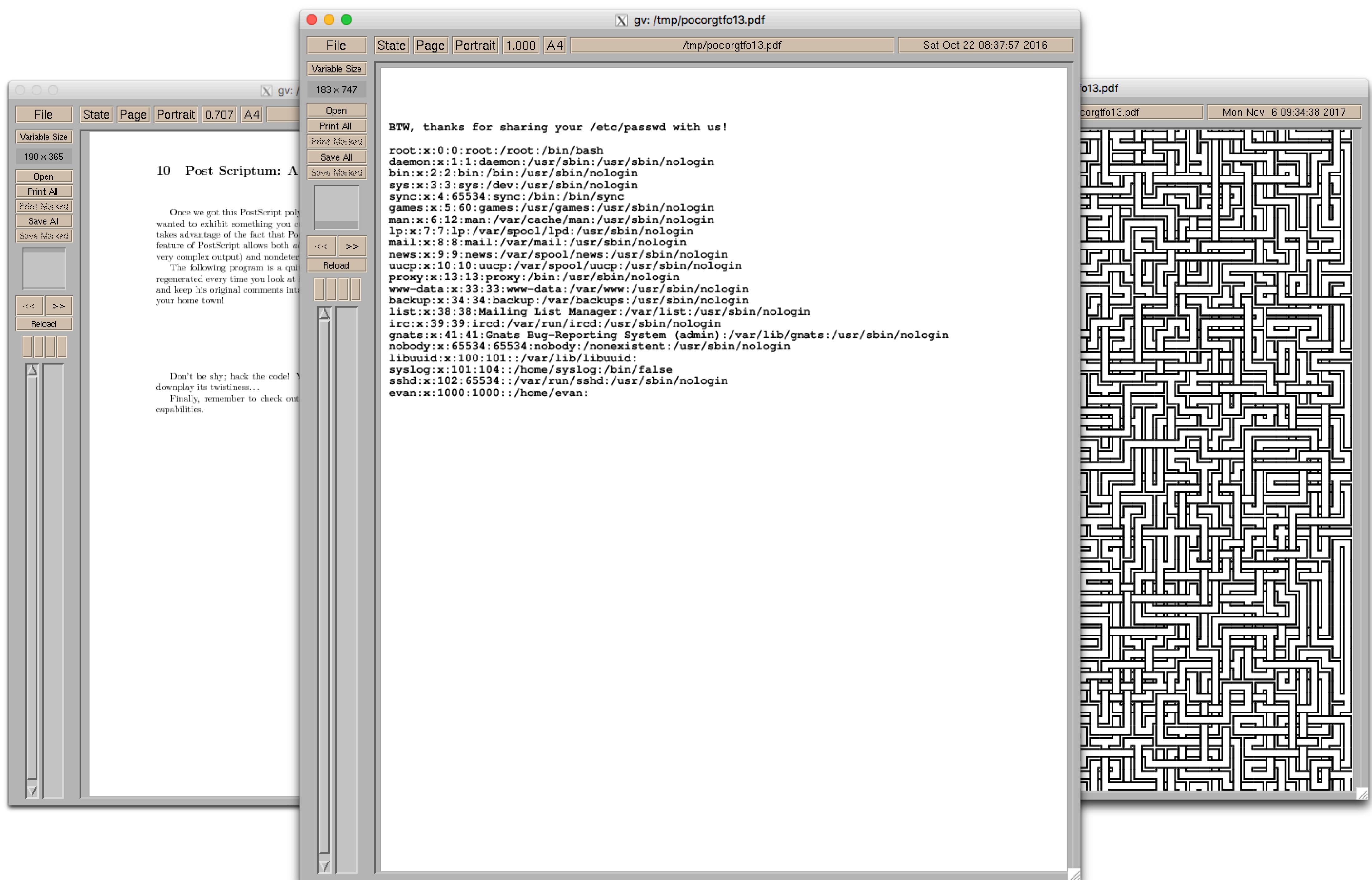
David Bau  
777 South Avenue  
Weston, MA 02193  
USA  
bau@cs.cornell.edu

Don't be shy; hack the code! You can tune this puzzle into an easier one if you disable tunneling or downplay its twistiness...

Finally, remember to check out the page after the puzzle for a closing demonstration of PostScript's capabilities.

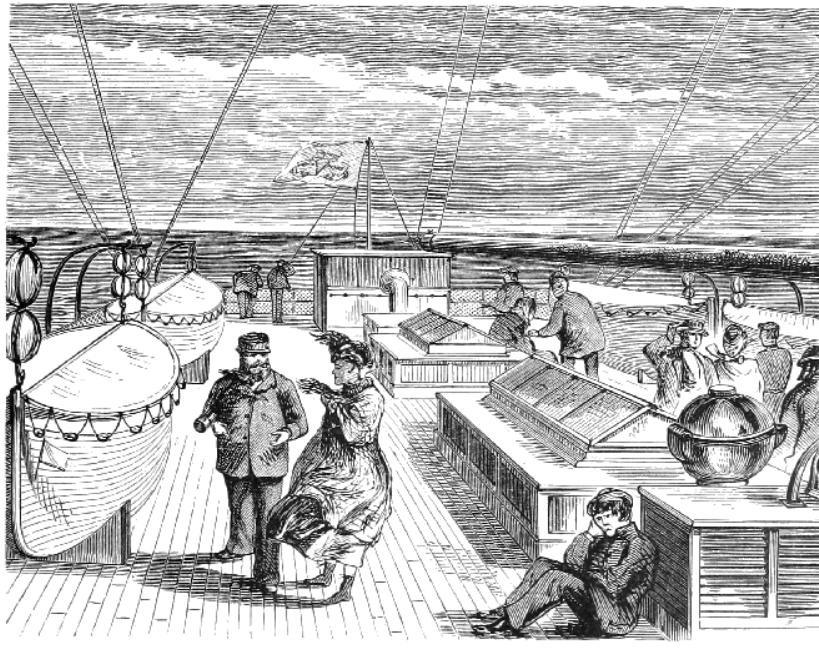
PostScript





# HTTP Quine

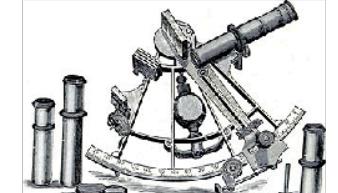
PoC||GTFO



IN A FIT OF STUBBORN OPTIMISM,  
PASTOR MANUL LAPHROAIG  
AND HIS CLEVER CREW  
SET SAIL TOWARD  
WELCOMING SHORES OF  
THE GREAT UNKNOWN!

11:1 Please Stand and Be Seated      11:6 Phrasebook for ARM Cortex M  
11:2 In Praise of Junk Hacking      11:7 Ghetto CFI for x86  
11:3 Emulating Star Wars on a Vector Display      11:8 Tourist's Guide to the MSP430  
11:4 Tron in 512 Bytes      11:9 This PDF is a Webserver  
11:5 Defeating the E7 Protection      11:10 In Memoriam: Ben "bushing" Byer  
Heidelberg, Baden-Württemberg

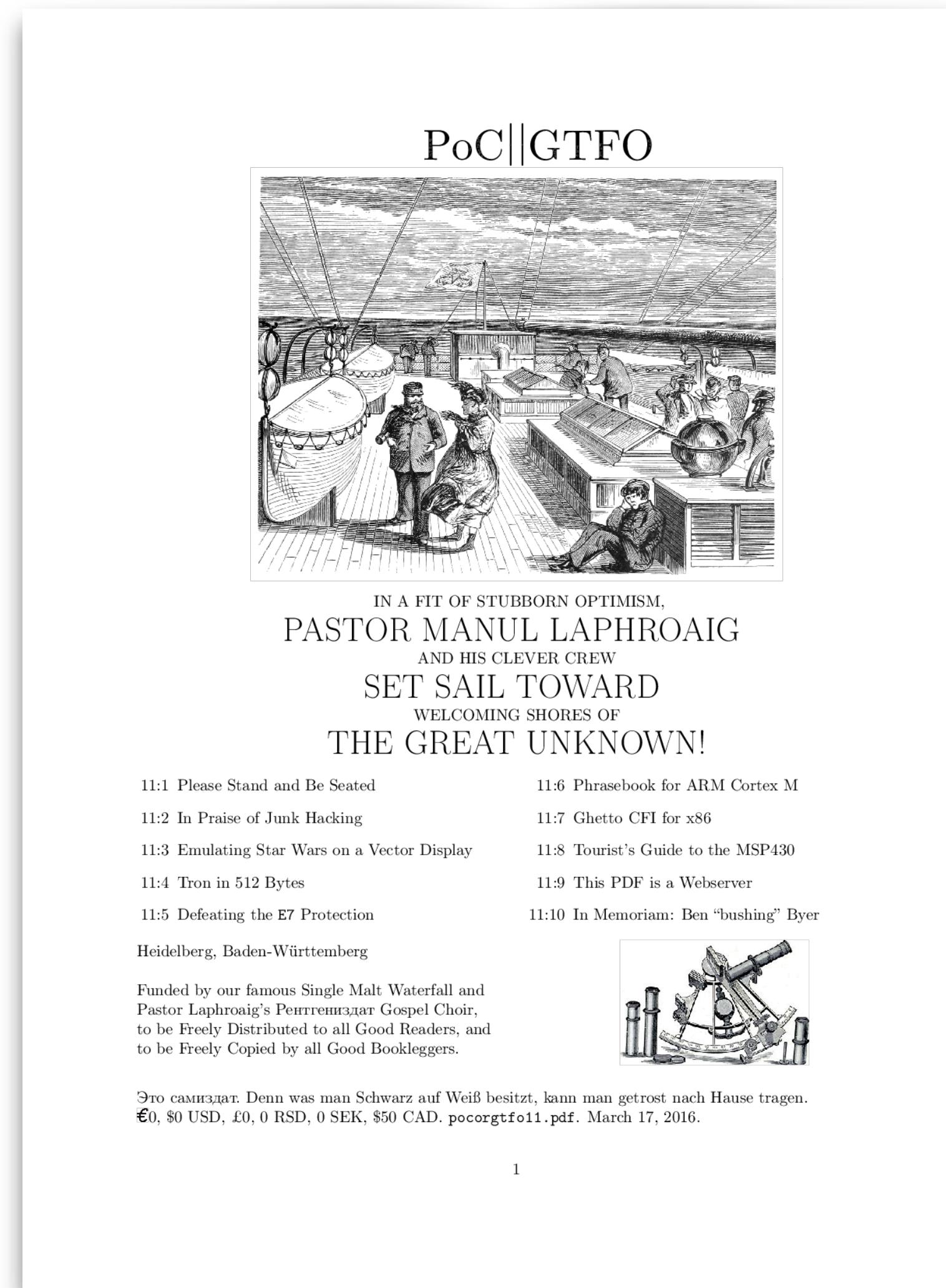
Funded by our famous Single Malt Waterfall and  
Pastor Laphroaig's Рентгениздат Gospel Choir,  
to be Freely Distributed to all Good Readers, and  
to be Freely Copied by all Good Bookleggers.



Это самиздат. Denn was man Schwarz auf Weiß besitzt, kann man getrost nach Hause tragen.  
€0, \$0 USD, £0, 0 RSD, 0 SEK, \$50 CAD. pocorgtfo11.pdf. March 17, 2016.

1

# HTTP Quine



```
$ ruby pocorgtfo11.pdf  
Listening for connections on port 8080.  
To listen on a different port,  
re-run with the desired port as a command-line argument.
```

The image shows a web browser window with the URL "localhost" in the address bar. The page content is identical to the PDF cover above, featuring the title "International Journal of PoC||GTFO Issue 0x11" and the same descriptive text. It also includes a "Click here to download the PDF!" button and a list of 10 articles. Below the article list, a note explains the polyglot nature of the page. On the right side, there's a section titled "Feelies" with a list of links to various files like "index.txt", "issues.txt", and "batterfirmware.pdf".

IN A FIT OF STUBBORN OPTIMISM,  
**PASTOR MANUL LAPHROAIG**  
AND HIS CLEVER CREW  
**SET SAIL TOWARD**  
WELCOMING SHORES OF  
**THE GREAT UNKNOWN!**

[Click here to download the PDF!](#)

11.1 Please Stand and Be Seated      11.6 Phrasebook for ARM Cortex M  
11.2 In Praise of Junk Hacking      11.7 Ghetto CFI for x86  
11.3 Emulating Star Wars on a Vector Display      11.8 Tourist's Guide to the MSP430  
11.4 Tron in 512 Bytes      11.9 This PDF is a Webserver  
11.5 Defeating the E7 Protection      11.10 In Memoriam: Ben "bushing" Byer

This is an HTML/Ruby/PDF/ZIP polyglot. When interpreted by Ruby, it acts as a web server that serves this page. If the URL requested of the webserver matches a path inside the ZIP, that file is served. If loaded directly in a web browser, the file will render as this webpage, too, however, the link to the PDF download is hidden.

**Feelies**

- [index.txt](#) — a text version of this feelies index
- [issues.txt](#) — about all issues of PoC||GTFO
- [issues.bib](#) — a BibTeX file containing references for all issues of PoC||GTFO
- [batteryfirmware.pdf](#) — Battery Firmware Hacking - Charlie Miller
- [vst.tar.bz2](#) — v.st vector board sources
- [vectormame.diff](#) — diff for Mame (see Star Wars article)
- [sluu225.pdf](#) — bq803xx ROM API v 3.0
- [favicon.png](#) — the icon for the website
- [bq20z80.py](#) — BQ20Z80 IDA Processor module

# HTTP Quine

The screenshot shows a web browser window titled "localhost" displaying the homepage of the International Journal of PoC||GTFO Issue 0x11. The page features a large title "International Journal of PoC||GTFO Issue 0x11" and a subtitle with a nautical theme: "IN A FIT OF STUBBORN OPTIMISM, PASTOR MANUL LAPHROAIG AND HIS CLEVER CREW SET SAIL TOWARD WELCOMING SHORES OF THE GREAT UNKNOWN!". Below this is a button labeled "Click here to download the PDF!". A sidebar on the left contains a small illustration of a ship and some text about the journal's history and funding. A sidebar on the right contains a green box with the word "argument.". The bottom of the page includes a "Feelies" section with links to various files and a note about the ZIP polyglot nature of the page.

## International Journal of PoC||GTFO Issue 0x11

IN A FIT OF STUBBORN OPTIMISM,  
PASTOR MANUL LAPHROAIG  
AND HIS CLEVER CREW  
SET SAIL TOWARD  
WELCOMING SHORES OF  
THE GREAT UNKNOWN!

[Click here to download the PDF!](#)

**11.1** Please Stand and Be Seated      **11.6** Phrasebook for ARM Cortex M  
**11.2** In Praise of Junk Hacking      **11.7** Ghetto CFI for x86  
**11.3** Emulating Star Wars on a Vector Display      **11.8** Tourist's Guide to the MSP430  
**11.4** Tron in 512 Bytes      **11.9** This PDF is a Webserver  
**11.5** Defeating the E7 Protection      **11.10** In Memoriam: Ben "bushing" Byer

This is an HTML/Ruby/PDF/ZIP polyglot. When interpreted by Ruby, it acts as a web server that serves this page. If the URL requested of the webserver matches a path inside the ZIP, that file is served. If loaded directly in a web browser, the file will render as this webpage, too, however, the link to the PDF download is hidden.

---

**Feelies**

- [index.txt — a text version of this feelies index](#)
- [issues.txt — about all issues of PoC||GTFO](#)
- [issues.bib — a BibTeX file containing references for all issues of PoC||GTFO](#)
- [batteryfirmware.pdf — Battery Firmware Hacking - Charlie Miller](#)
- [vst.tar.bz2 — v.st vector board sources](#)
- [vectormame.diff — diff for Mame \(see Star Wars article\)](#)
- [sluu225.pdf — bq803xx ROM API v 3.0](#)
- [favicon.png — the icon for the website](#)
- [bq20z80.py — BQ20Z80 IDA Processor module](#)

**Feelies automagically parsed from the ZIP!**

# But Wait, There's More!

```
$ ln -s pocorgtfo11.pdf pocorgtfo11.html
```

# But Wait, There's More!

A screenshot of a web browser window displaying the homepage of the International Journal of PoC||GTFO Issue 0x11. The browser has a light gray header with standard OS X-style buttons and a search bar containing the URL "pocorgtfo11.html". The main content area features a large, bold title "International Journal of PoC||GTFO Issue 0x11". Below it is a subtitle in a smaller, bold font: "IN A FIT OF STUBBORN OPTIMISM,  
PASTOR MANUL LAPHROAIG  
AND HIS CLEVER CREW  
SET SAIL TOWARD  
WELCOMING SHORES OF  
THE GREAT UNKNOWN!". The page lists ten articles under two columns: "11.1 Please Stand and Be Seated", "11.2 In Praise of Junk Hacking", "11.3 Emulating Star Wars on a Vector Display", "11.4 Tron in 512 Bytes", "11.5 Defeating the E7 Protection", "11.6 Phrasebook for ARM Cortex M", "11.7 Ghetto CFI for x86", "11.8 Tourist's Guide to the MSP430", "11.9 This PDF is a Webserver", and "11.10 In Memoriam: Ben ‘bushing’ Byer". At the bottom, there is a note about the polyglot nature of the page, information about the location (Heidelberg, Baden-Württemberg), funding details, and a quote in Russian.

**International Journal of PoC||GTFO Issue 0x11**

IN A FIT OF STUBBORN OPTIMISM,  
**PASTOR MANUL LAPHROAIG**  
AND HIS CLEVER CREW  
**SET SAIL TOWARD**  
WELCOMING SHORES OF  
**THE GREAT UNKNOWN!**

**11.1** Please Stand and Be Seated      **11.6** Phrasebook for ARM Cortex M  
**11.2** In Praise of Junk Hacking      **11.7** Ghetto CFI for x86  
**11.3** Emulating Star Wars on a Vector Display      **11.8** Tourist's Guide to the MSP430  
**11.4** Tron in 512 Bytes      **11.9** This PDF is a Webserver  
**11.5** Defeating the E7 Protection      **11.10** In Memoriam: Ben “bushing” Byer

This is an HTML/Ruby/PDF/ZIP polyglot. When interpreted by Ruby, it acts as a web server that serves this page. If the URL requested of the webserver matches a path inside the ZIP, that file is served. If loaded directly in a web browser, the file will render as this webpage, too, however, the link to the PDF download is hidden.

---

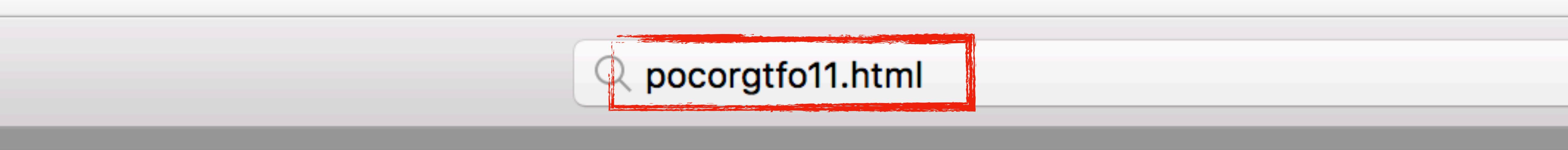
Heidelberg, Baden-Württemberg

Funded by our famous Single Malt Waterfall and  
Pastor Laphroaig’s Рентгениздат Gospel Choir,  
to be Freely Distributed to all Good Readers, and  
to be Freely Copied by all Good Bookleggers.

Это самиздат. Denn was man Schwarz auf Weiß besitzt, kann man getrost nach Hause tragen.

€0, \$0 USD, £0, 0 RSD, 0 SEK, \$50 CAD. March 13, 2016.

# But Wait, There's More!



# International Journal of Po

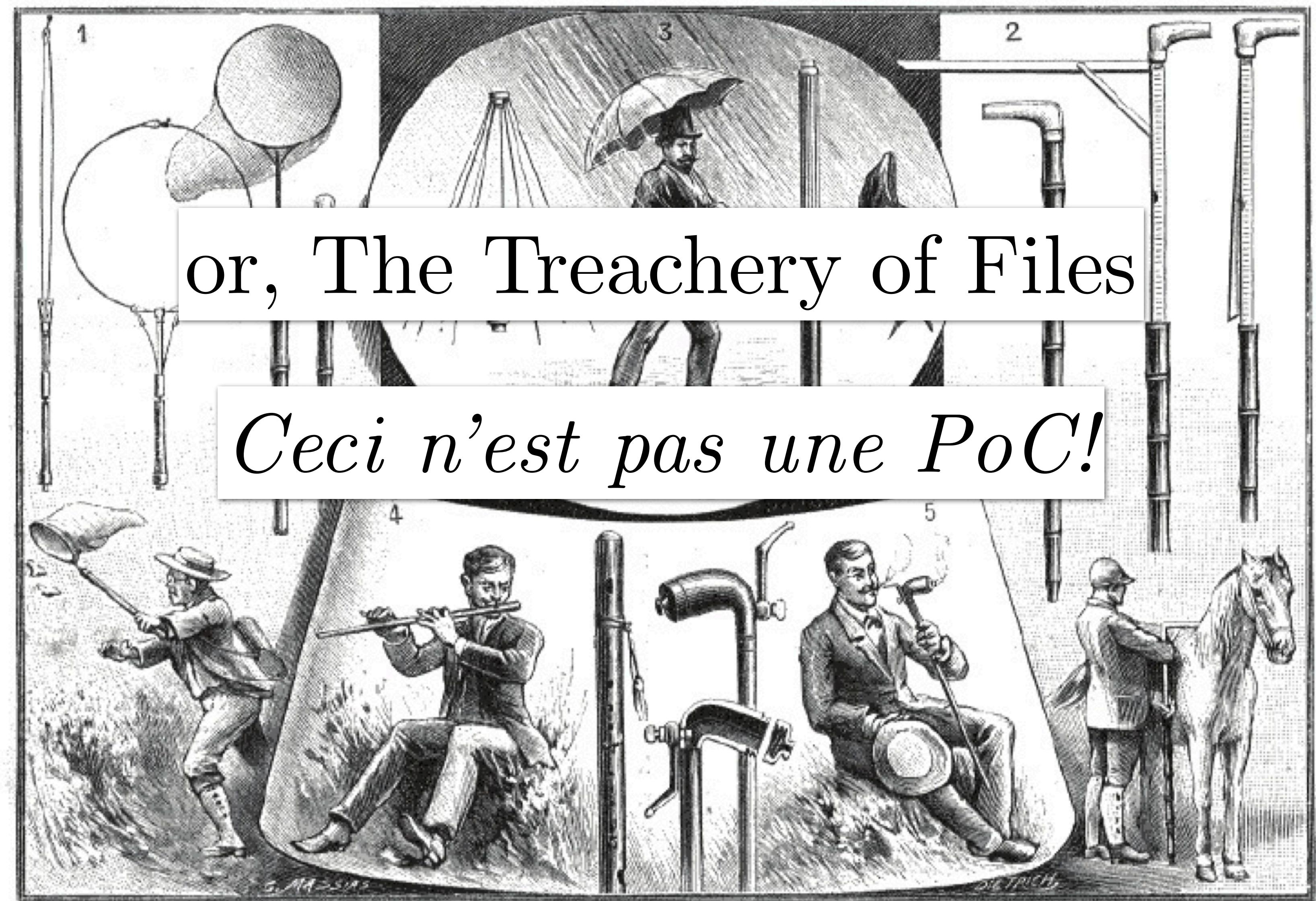
IN A FIT OF STUBBORN OPTI  
PASTOR MANUEL JAP

Which is also a Ruby script

In which an HTML page is also a PDF

Which is also a ZIP

Which is an HTTP Quine



Utilisation de la canne. — 1. Canne-filet à papillons. — 2. Canne à toiser les chevaux. —  
3. Canne-parapluie. — 4. Canne musicale. — 5. Ceci n'est pas une pipe.

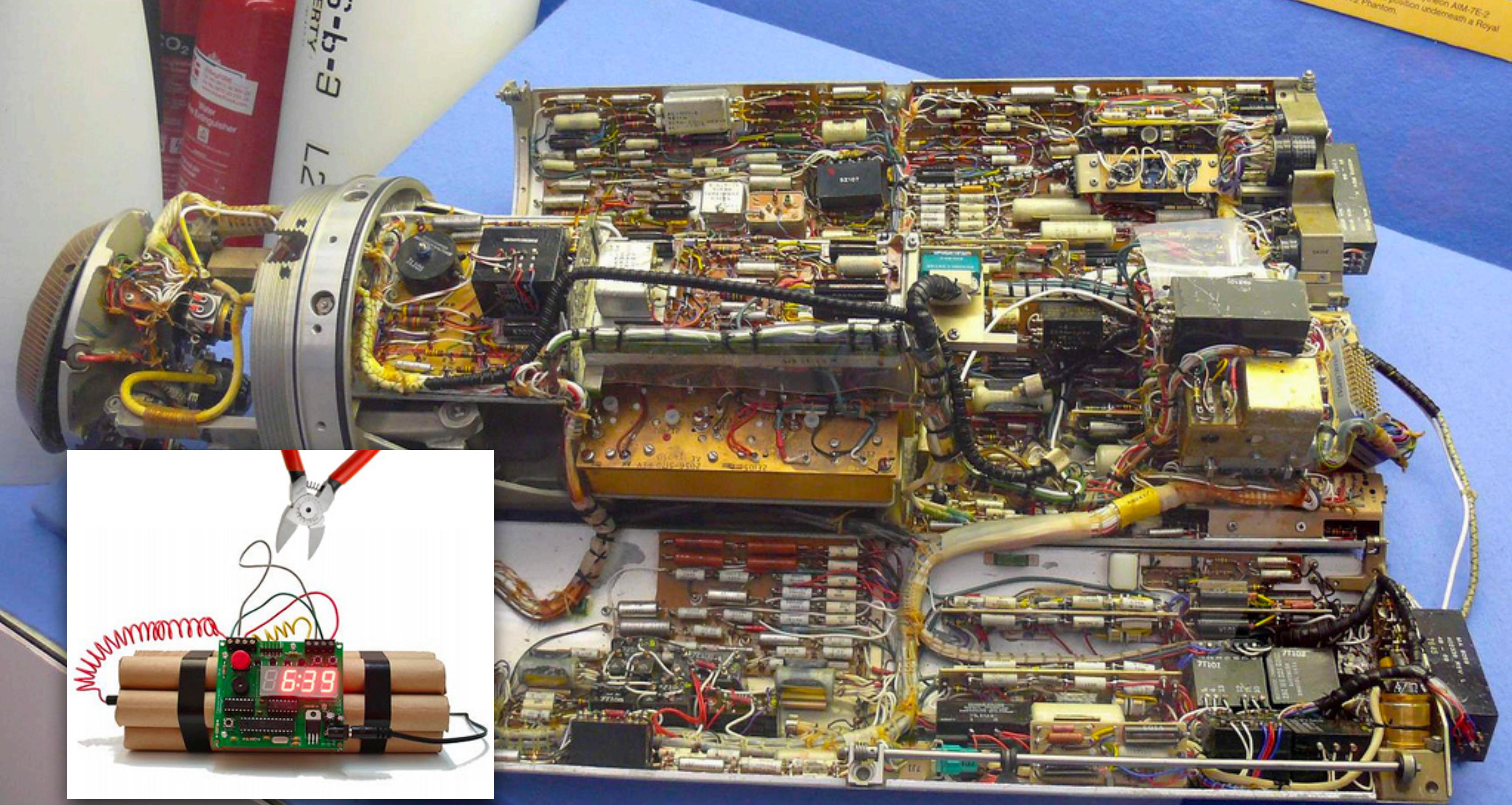


3,000,000 Square Feet  
(almost 2x the size of the  
Tesla Gigafactory)

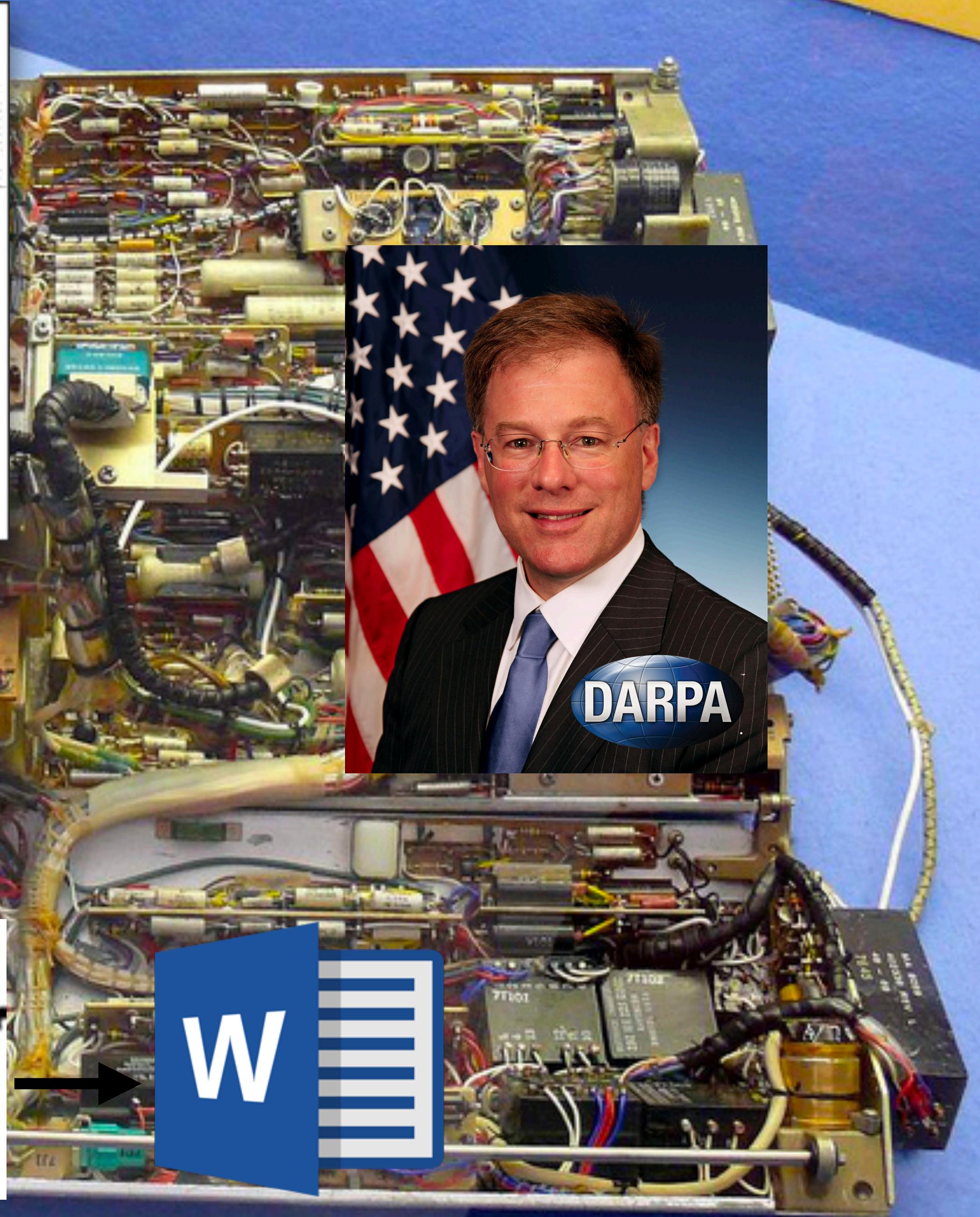


1949–2014





# Preserved external research. (blog ⇒ PDF)



# PDF/Git Polyglot

```
$ git bundle
```

# bundle.c

```
argv_array_pushl(&pack_objects.args,
    "pack-objects", "--all-progress-implied",
    "--compression=0",
    "--stdout", "--thin", "--delta-base-offset",
NULL);
```

```
$ export PATH=/path/to/patched/git:$PATH
$ git init
$ git add article.pdf
$ git commit article.pdf -m "added"
$ git bundle create PDFGitPolyglot.pdf -all
```

This PDF is a Git Repository  
Containing its Own L<sup>A</sup>T<sub>E</sub>X Source  
and a Copy of Itself

Evan Sultanik

April 11, 2017

Have you ever heard of the `git bundle` command? I hadn't. It bundles a set of Git objects—potentially even an entire repository—into a single file. Git allows you to treat that file as if it were a standard Git database, so you can do things like clone a repo directly from it. Its purpose is to easily sneakernet pushes or even whole repositories across air gaps.

---

Neighbors, it's possible to create a PDF that is also a Git repository.

# PDF/Git Polyglot

This PDF is a Git Repository  
Containing its Own L<sup>A</sup>T<sub>E</sub>X Source  
and a Copy of Itself

Evan Sultanik  
April 11, 2017

Have you ever heard of the `git bundle` command? I hadn't. It bundles a set of Git objects—potentially even an entire repository—into a single file. Git allows you to treat that file as if it were a standard Git database, so you can do things like clone a repo directly from it. Its purpose is to easily sneakernet pushes or even whole repositories across air gaps.

```
-----  
Neighbors, it's possible to create a PDF that is also a Git repository.  
-----  
$ git clone PDFGitPolyglot.pdf foo  
Cloning into 'foo'...  
Receiving objects: 100% (174/174), 103.48 KiB | 0 bytes/s, done.  
Resolving deltas: 100% (100/100), done.  
$ ls  
PDFGitPolyglot.pdf PDFGitPolyglot.tex
```

## 1 The Git Bundle File Format

The file format for Git bundles doesn't appear to be formally specified anywhere, however, inspecting `bundle.c` reveals that it's relatively straightforward:

```
-----  
# v2 git bundle file  
-----  
[Git Bundle Signature]  
-----  
3aa340a2e3d125b6703e5c9b1f1de2054a9c0c refs/heads/master  
3aa340a2e3d125b6703e5c9b1f1de2054a9c0c refs/remotes/origin/master  
4146cfe2fe9249fc146232832587efef197ef5d2d refs/attach  
babdd4a4735ee1f64b7023be3545860d0db0bae250a HEAD  
-----  
[PACK...]  
-----  
[Git Packfile]
```

Git has another custom format called a *Packfile* that it uses to compress the objects in its database, as well as to reduce network bandwidth when pushing

```
$ git clone PDFGitPolyglot.pdf foo  
Cloning into 'foo'...  
Receiving objects: 100% (174/174), 103.48 KiB | 0 bytes/s, done.  
Resolving deltas: 100% (100/100), done.  
$ cd foo  
$ ls  
PDFGitPolyglot.pdf PDFGitPolyglot.tex
```

# Act II

Introducing PolyFile and PolyTracker!

# Act II

## Introducing PolyFile and PolyTracker!

**Automated Lexical Annotation  
and Navigation of Parsers**

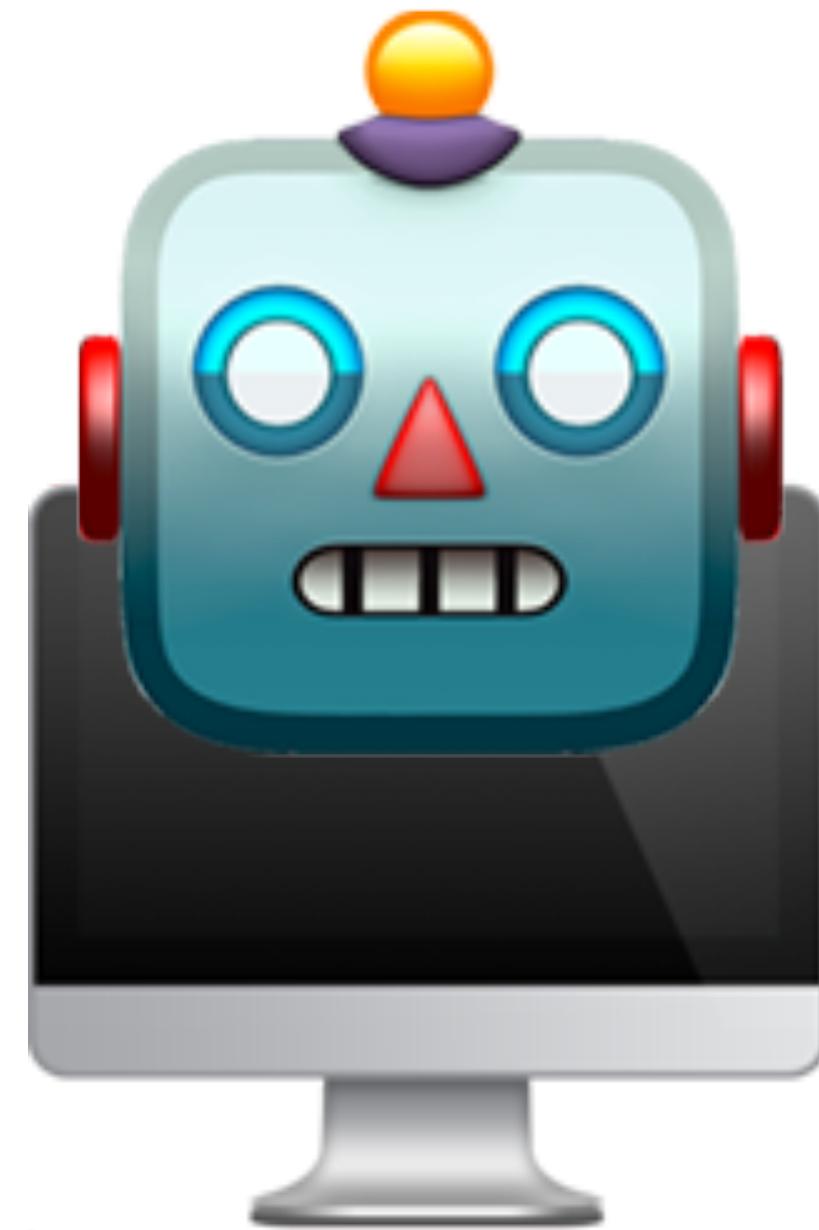
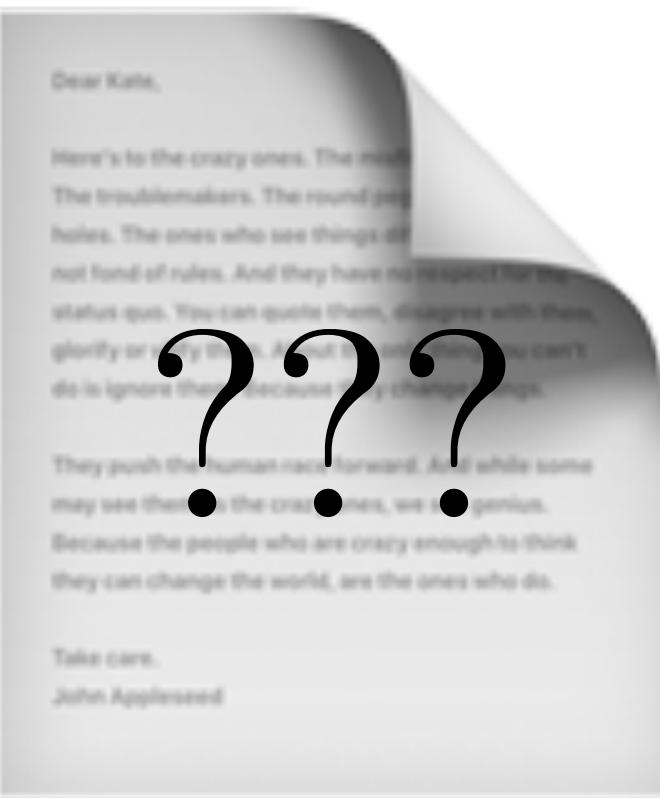
# Act II

## Introducing PolyFile and PolyTracker!

The ALAN Parsers Project  
~~Automated Lexical Annotation  
and Navigation of Parsers~~

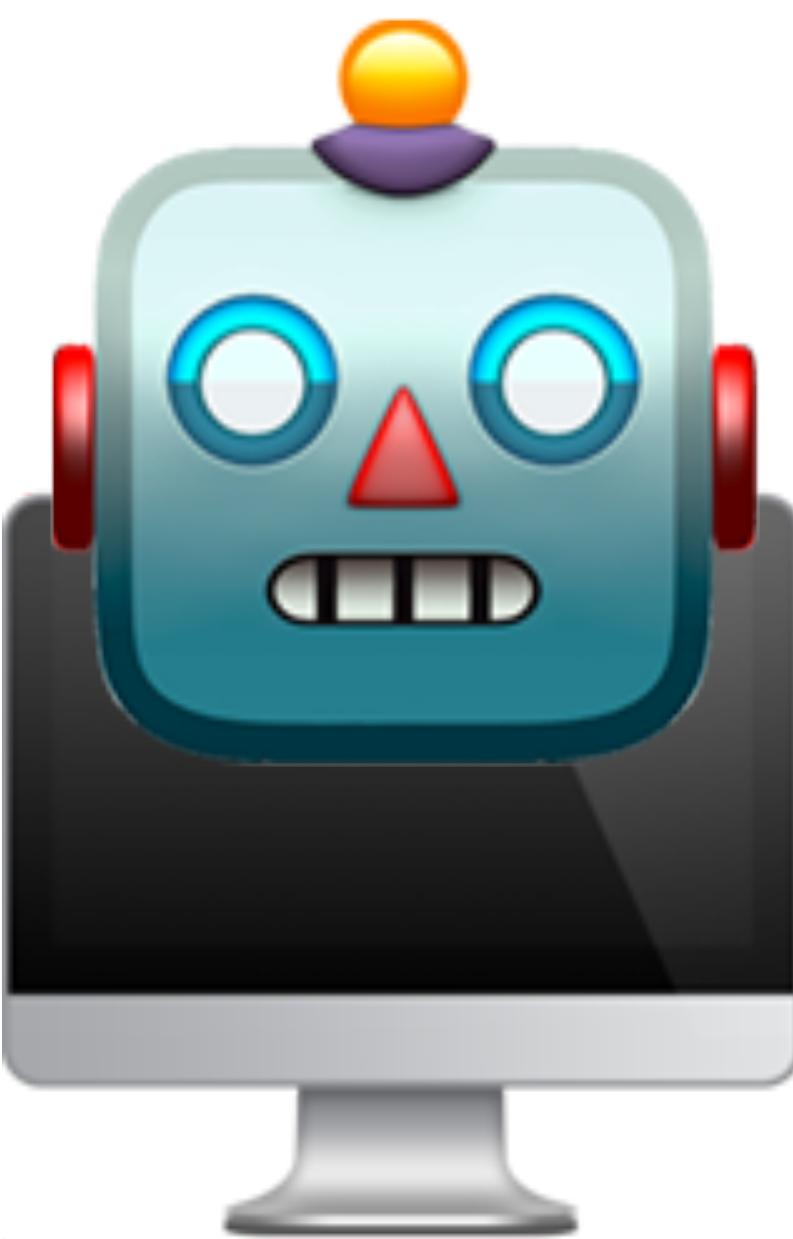
# High Level Goals

Create semantic map of the functions in a parser



# High Level Goals

Create semantic map of the functions in a parser



parser\_function1

↳byte 0, 10, 50

Object Stream

parser\_function2

↳byte 10, 74

Xref

parser\_function3

↳byte 20

JFIF

# High Level Goals

Create semantic map of the functions in a parser

parser\_function1

↳byte 0, 10, 50

Object Stream

parser\_function2

↳byte 10, 74

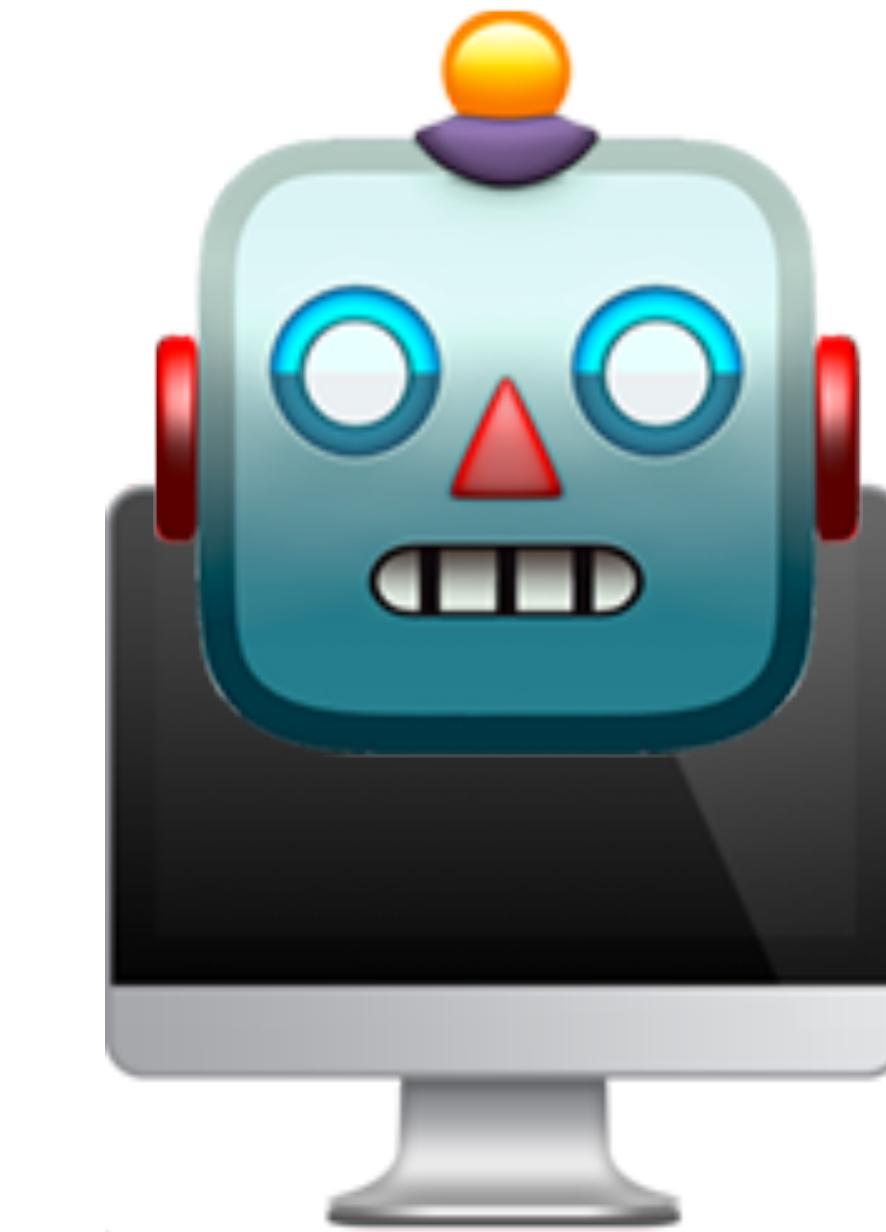
Xref

parser\_function3

↳byte 20

JFIF

**Ultimate Goal:** Automatically extract a minimal grammar specifying the files accepted by a parser



**Hypothesis:** The majority of the potential for maliciousness and schizophrenia will exist in the symmetric difference of the grammars accepted by a format's parser implementations

# Related Work

- **Ange's SBuD**
  - Limited parser support
- **Kaitai Struct**
  - No support for PDF
  - Difficult to specify context sensitive languages
- **Valgrind / TaintGrind**
  - Intractably slow; hours of computation to parse a minimal PDF
- **AFL Analyze**
  - Does not associate input byte sequences with the program trace
- **LLVM Dataflow Sanitizer**
  - Limited to just a few thousand taint labels
  - Can only track at most one input byte at a time

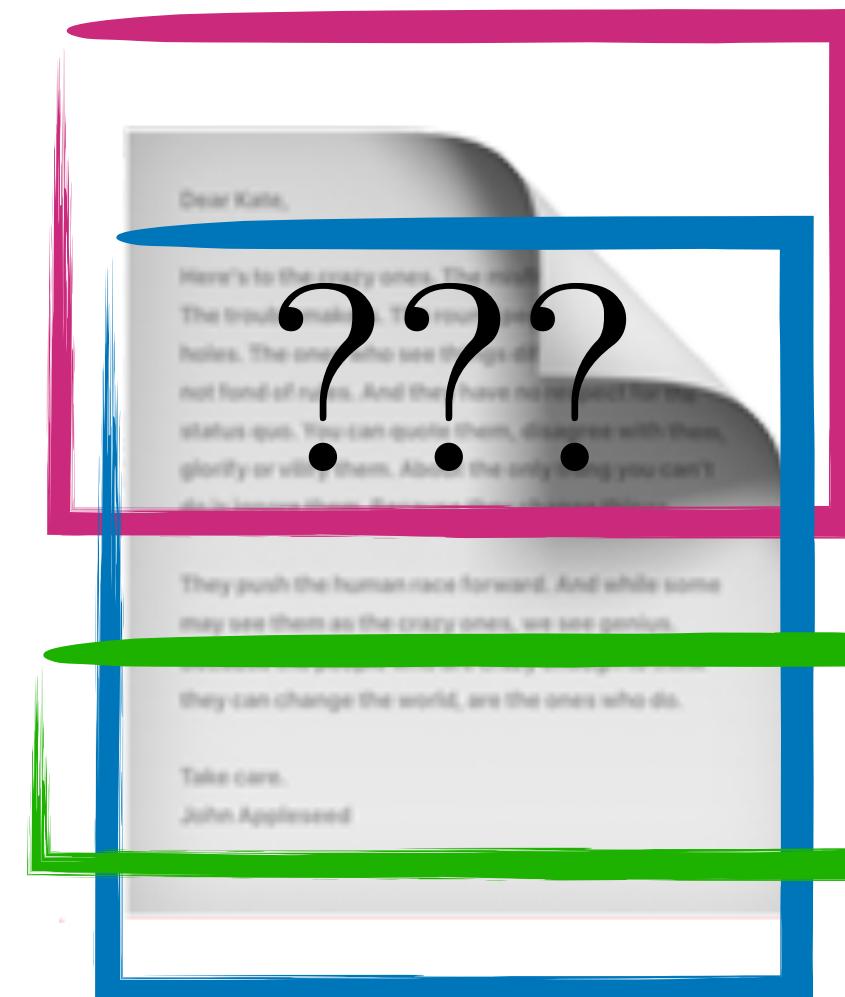


Valgrind



# Approach: Semantic Labeling

## Polyglot-Aware File Identification



iNES ROM  
PDF  
ZIP

Modify parsers for best effort

Instrument to track input byte offsets

Label regions of the input

Produce ground truth

## Hierarchical Labeling

iNES [0x0→0x12220]  
↳ Magic [0x0→0x3]  
Header [0x4→0xF]  
:  
PRG [0xC210→0x1020F]  
CHR [0x10210→0x12220]  
PDF [0x10→0x2EF72F]  
↳ Magic [0x10→0x1E]  
Object 1.0 [0x1F→0x12221]  
↳ Dictionary [0x2A→0x3E]  
Stream [0x3F→0x12219]  
↳ JFIF Image [0x46→0x1220F]  
↳ JPEG Segment [...]  
↳ Magic [...]  
Marker [...]

# Approach: Parser Instrumentation

## LLVM

Operate on LLVM/IR

Can work with all open source parsers

Eventually support closed-source binaries by lifting to LLVM (*e.g.*, with McSEMA or Remill)

## Instrumentation

Shadow memory inspired by the Data Flow Sanitizer (dfsan)

Negligible CPU overhead

$O(n)$  memory overhead, where n is the number of instructions executed by the parser

## Taint Tracking

Novel datastructure for efficiently storing taint labels

dfsan status quo:

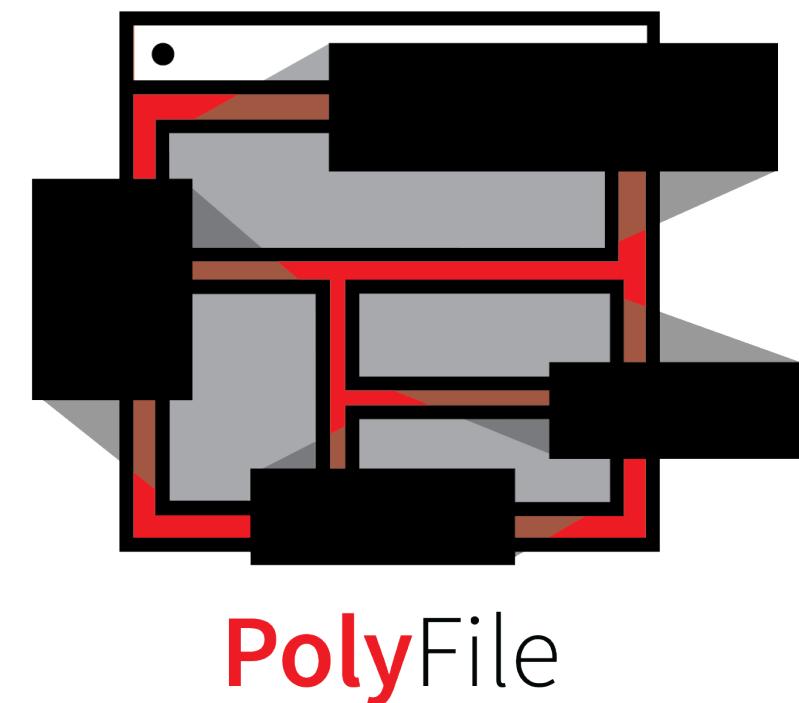
$\Theta(1)$  lookups

$\Theta(n^2)$  storage

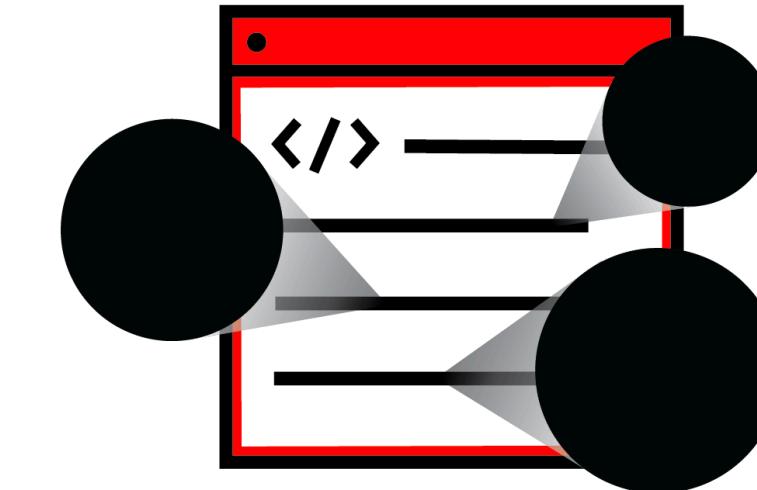
TAPP:

$O(\log n)$  lookups

$O(n)$  storage



# Current Tooling



PolyTracker



Novel file matching algorithm seeded with over 10,000 file definitions from TrID

Shadow memory inspired by the Data Flow Sanitizer (dfsan) and Angora



Kaitai Struct file definition parser for automatically labeling based upon KSY definitions (*e.g.*, JPEG, ZIP, WASM)

Negligible CPU overhead



$O(n)$  memory overhead, where n is the number of instructions executed by the parser



Instrumented Didier Stevens' PDF parser, resilient to malformations and schizophrenia

Currently supports C, C++ (experimental), and any code entirely representable in LLVM/IR.

Interactive file explorer

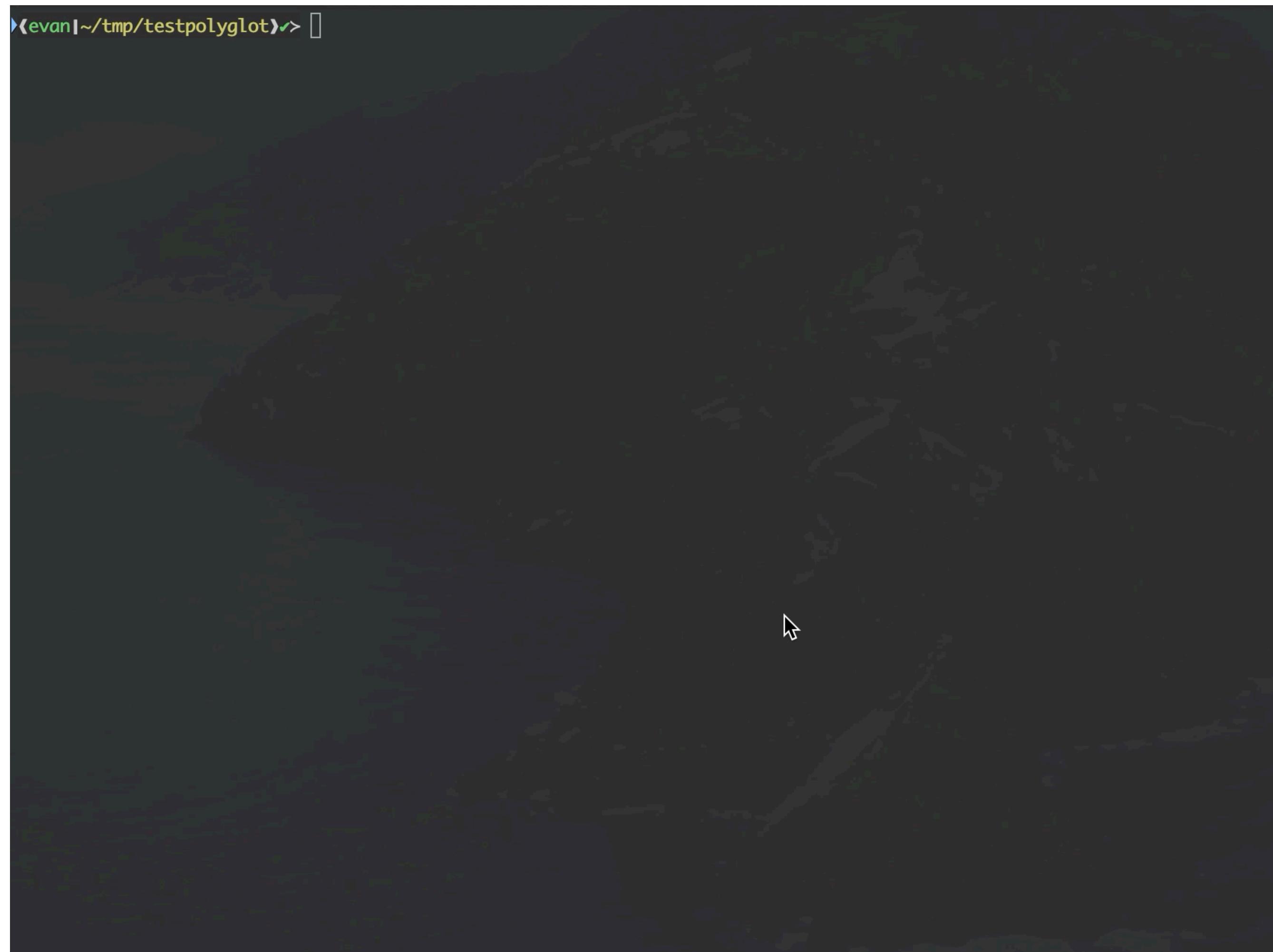
<https://github.com/trailofbits/polyfile>

<https://github.com/trailofbits/polytracker>

# Demo

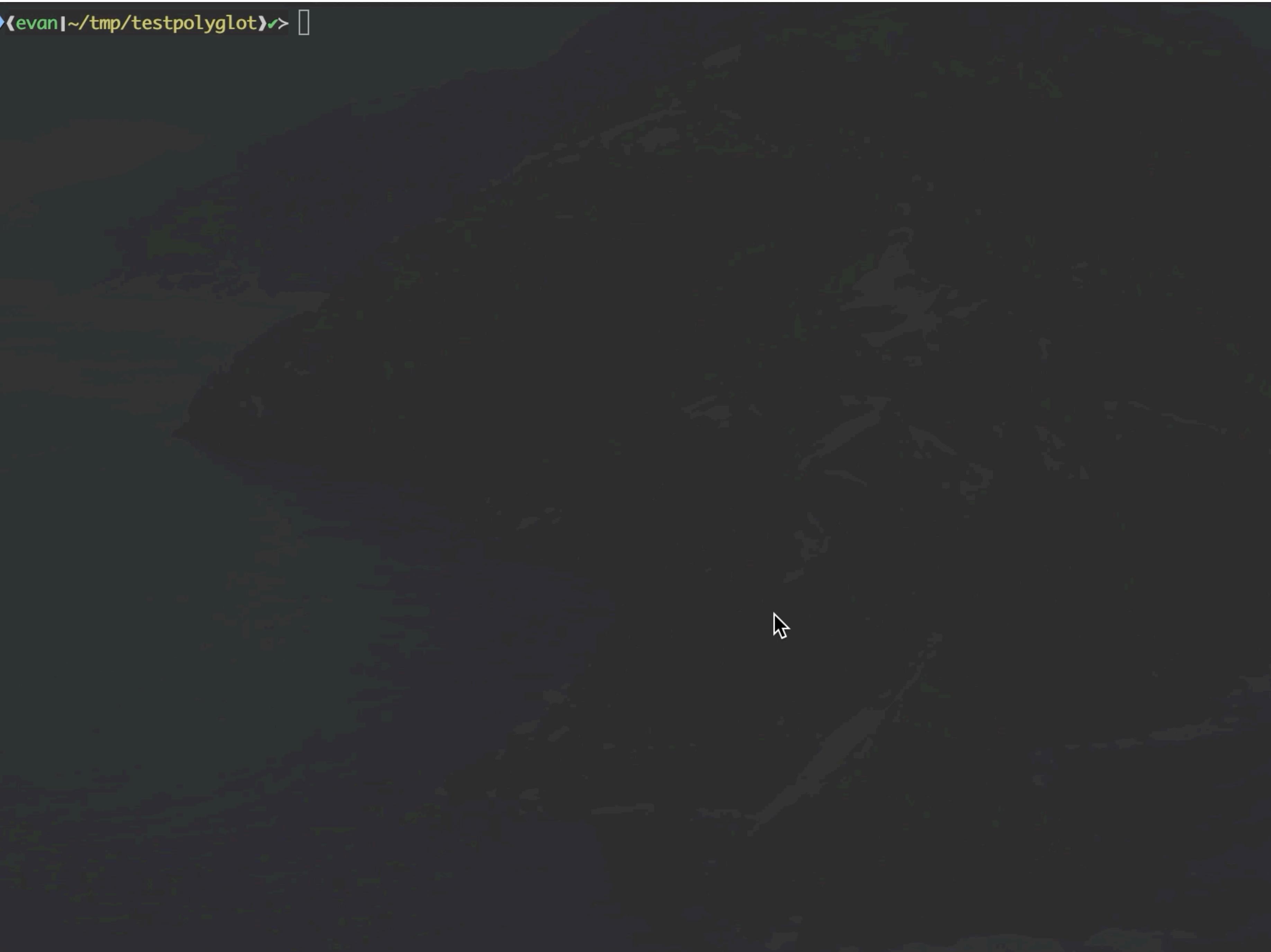
# Demo

# WASM/HTML Polyglot Detection



# WASM/HTML

## Polyglot Detection

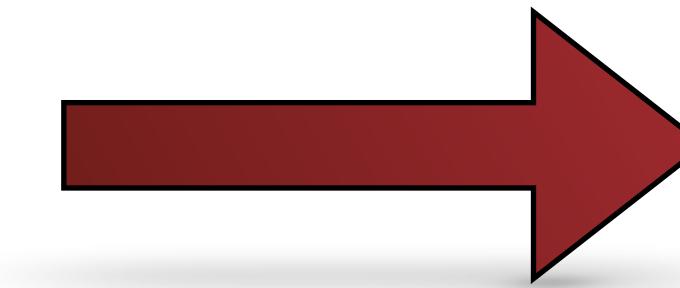


# PolyTracker Instrumentation

```
{  
    "dfs$ensure_solid_xref": [  
        2276587,  
        2276588  
    ],  
    "dfs$fmt_obj": [  
        2465223,  
        2465224,  
        2465225,  
        2465226,  
        2465227,  
        2465228,  
        2465240,  
        2465241,  
        2465242,  
        2465243,  
        2465244,  
        2465245,  
        2465246,  
        2465258,  
        2465259,  
        2465260,  
        2465261,  
        2465262  
    ]  
}
```

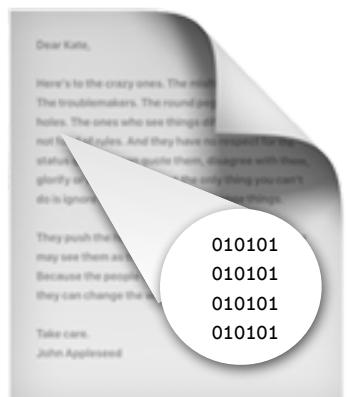
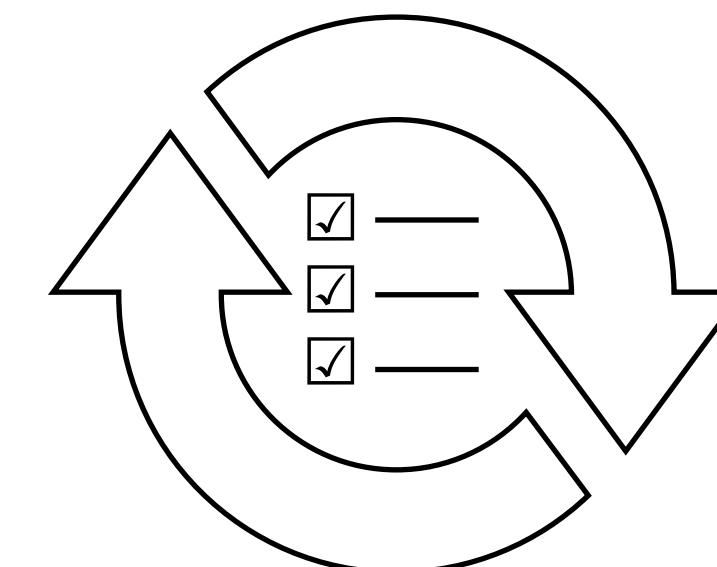
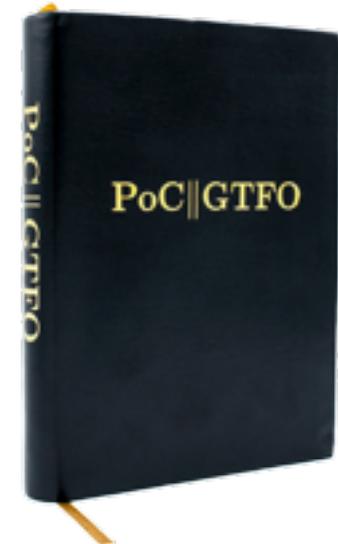
# PolyTracker Instrumentation

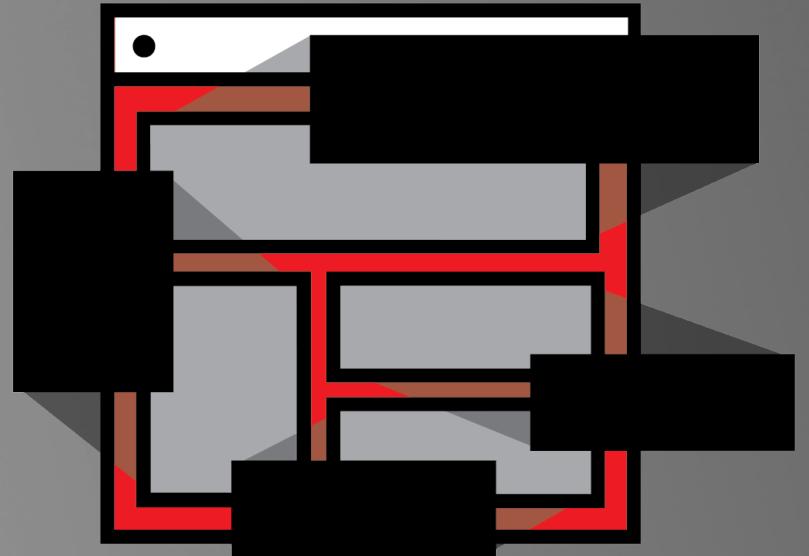
```
{  
    "dfs$ensure_solid_xref": [  
        2276587,  
        2276588  
    ],  
    "dfs$fmt_obj": [  
        2465223,  
        2465224,  
        2465225,  
        2465226,  
        2465227,  
        2465228,  
        2465240,  
        2465241,  
        2465242,  
        2465243,  
        2465244,  
        2465245,  
        2465246,  
        2465258,  
        2465259,  
        2465260,  
        2465261,  
        2465262  
    ]  
}
```



iNES [0x0→0x12220]  
↳ Magic [0x0→0x3]  
    Header [0x4→0xF]  
    :  
    PRG [0xC210→0x1020F]  
    CHR [0x10210→0x12220]  
    PDF [0x10→0x2EF72F]  
↳ Magic [0x10→0x1E]  
    Object 1.0 [0x1F→0x12221]  
    ↳ Dictionary [0x2A→0x3E]  
        Stream [0x3F→0x12219]  
        ↳ JFIF Image [0x46→0x1220F]  
            ↳ JPEG Segment [...]  
            ↳ Magic [...]  
            Marker [...]  
    :  
    :

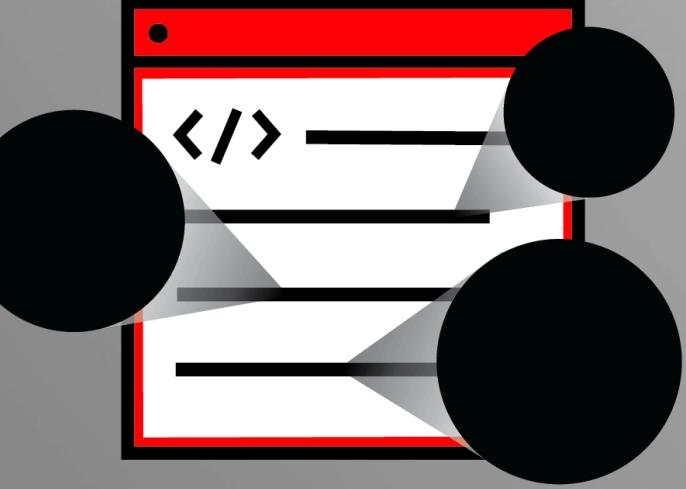
# Why Should You Care?



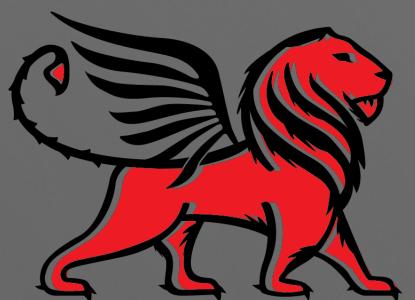
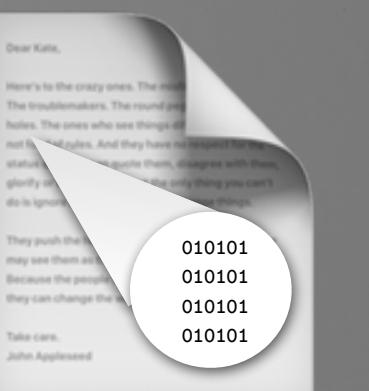
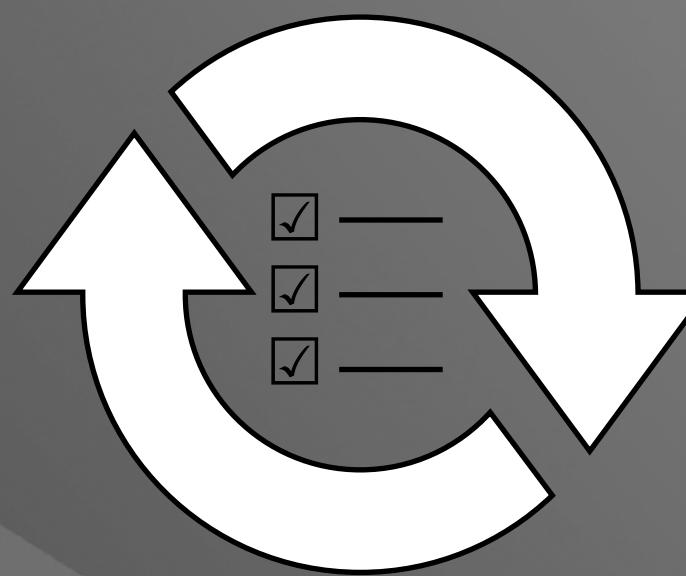
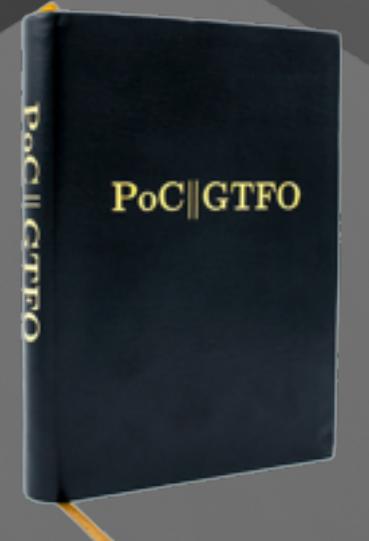


PolyFile

# Why Should You Care?



PolyTracker



# Acknowledgements

Ange Albertini

@angealbertini

Sergey Bratus

@sergeybratus

Travis Goodspeed

@travisgoodspeed

Philippe Teuwen

@doegox

Evan Teran

@evan\_teran

Jacob Torrey

@JacobTorrey

Ryan Speers

@rmspeers

*Et pl. al.*



Evan Sultanik

Thanks!

@ESultanik



# Thanks!



Carson Harmon

@reyeetengineer

Brad Larsen

@BradLarsen

Evan Sultanik

@ESultanik

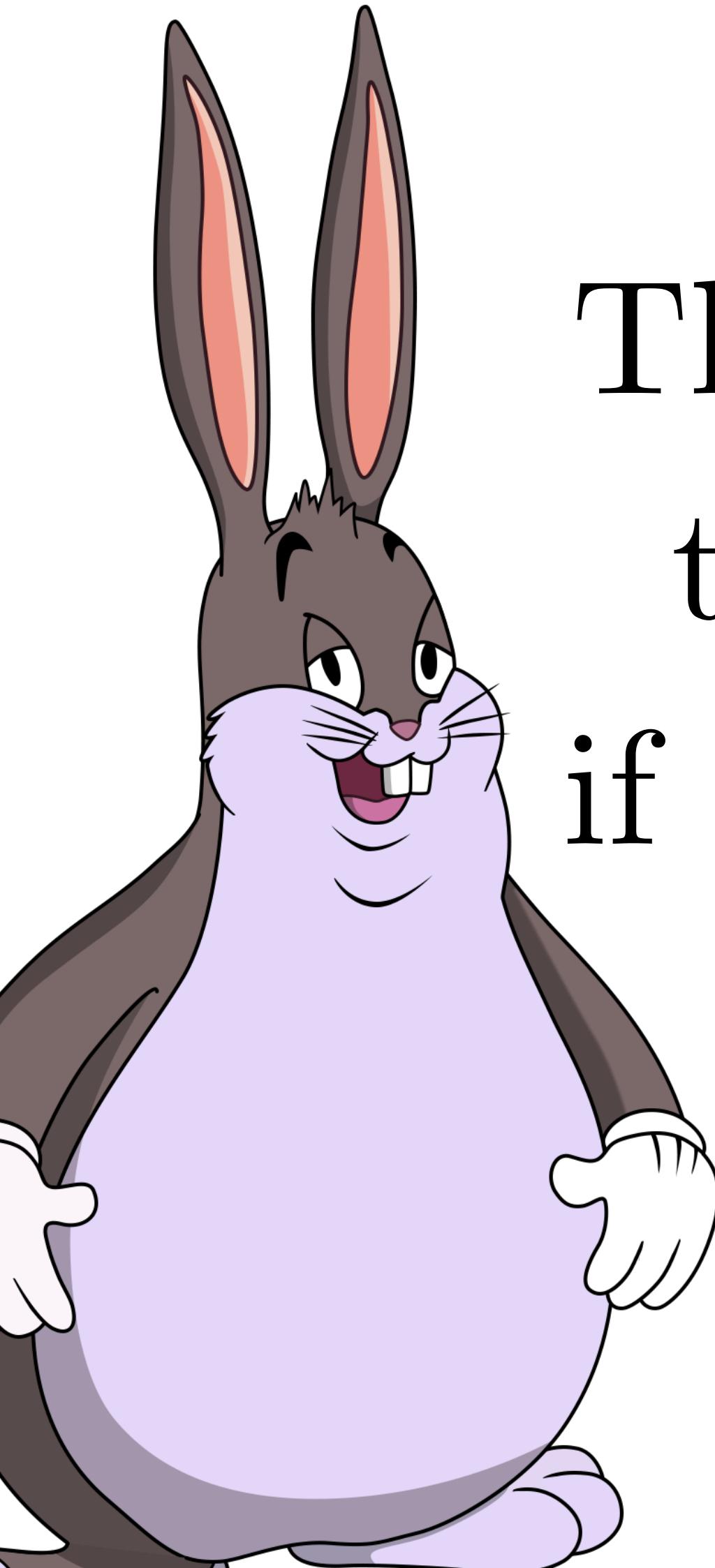


<https://github.com/trailofbits/polyfile>

<https://github.com/trailofbits/polytracker>

# Contact Info





This slide is here solely to prevent  
the slideshow from auto-closing  
if Evan clicks one too many times!