# Concepts and Tooling for Reverse Engineering

Florian Magin <fmagin@ernw.de>

# whoami

- Security Researcher at ERNW Research GmbH from Heidelberg, Germany

- Organizer of the Wizards of Dos CTF team from Darmstadt, Germany

- Reach me via:
  - Twitter: @0x464D
  - Email: fmagin@ernw.de

# Who we are

o Germany-based ERNW GmbH

o Blog: *www.insinuator.net*

o Conference: *www.troopers.de*

# Agenda

o How to extract meaning from a bunch of bytes

    o With a focus on what happens if the bytes contain executable code for a Linux like system

# Bunch of Bytes

o Find Patterns

    o The human brain is really good at that

o Throw some byte sequences into a search machine

o Contextualize

    o In most cases you know the rough context

o Just call 'file' or 'binwalk' on it

o Find a good enough parser

o If there is none, generate your own

    o More on that later

https://github.com/ReFirmLabs/binwalk

# First Steps

o  If it's in a typical executable format

    o  You are lucky

    o  Plenty of parsers and support

    o  Most information already available

o  If it's firmware

    o  Manual work required

# Firmware Information

o Determine the Architecture

    o Datasheet

    o Heuristics like grepping for function prologues/epilogues for various CPUs/CCs

o Determine Memory Layout

    o Datasheet

    o Memory Dump

# Basics

○ Using common Linux tooling and internals

# Executable Parsing

- ○ ELF: Executable and Linking Format
- ○ PE/MZ
- ○ Mach-0
- ○ Header contains all the information for the loader
  to setup the program
  - ○ Memory layout
  - ○ Entry point
  - ○ Dependencies
  - ○ etc
- ○ Relevant:
  - ○ External Library and Function names
  - ○ Symbol Table if available

# Executable Parsing: Tools

o  Many overlapping tools

o  'objdump' if you want to get a first look

o  Other tools will take care of this for you

```
/bin/ls:     file format elf64-x86-64
/bin/ls
architecture: i386:x86-64, flags 0x00000150:
HAS_SYMS, DYNAMIC, D_PAGED
start address 0x0000000000005000

Program Header:
    PHDR off    0x0000000000000040 vaddr 0x0000000000000040 paddr 0x0000000000000040 align 2**3
         filesz 0x00000000000001f8 memsz 0x00000000000001f8 flags r-x
  INTERP off    0x0000000000000238 vaddr 0x0000000000000238 paddr 0x0000000000000238 align 2**0
         filesz 0x000000000000001c memsz 0x000000000000001c flags r--
    LOAD off    0x0000000000000000 vaddr 0x0000000000000000 paddr 0x0000000000000000 align 2**21
         filesz 0x000000000001e050 memsz 0x000000000001e050 flags r-x
    LOAD off    0x000000000001f030 vaddr 0x000000000021f030 paddr 0x000000000021f030 align 2**21
         filesz 0x0000000000001238 memsz 0x0000000000002530 flags rw-
 DYNAMIC off    0x000000000001fa78 vaddr 0x000000000021fa78 paddr 0x000000000021fa78 align 2**3
         filesz 0x00000000000001c0 memsz 0x00000000000001c0 flags rw-
    NOTE off    0x0000000000000254 vaddr 0x0000000000000254 paddr 0x0000000000000254 align 2**2
         filesz 0x0000000000000044 memsz 0x0000000000000044 flags r--
EH_FRAME off    0x000000000001ac0c vaddr 0x000000000001ac0c paddr 0x000000000001ac0c align 2**2
         filesz 0x000000000000084c memsz 0x000000000000084c flags r--
   STACK off    0x0000000000000000 vaddr 0x0000000000000000 paddr 0x0000000000000000 align 2**4
         filesz 0x0000000000000000 memsz 0x0000000000000000 flags rw-
   RELRO off    0x000000000001f030 vaddr 0x000000000021f030 paddr 0x000000000021f030 align 2**0
         filesz 0x0000000000000fd0 memsz 0x0000000000000fd0 flags r--

Dynamic Section:
  NEEDED               libcap.so.2
  NEEDED               libc.so.6
  INIT                 0x00000000000035a8
  FINI                 0x0000000000015dfc
  INIT_ARRAY           0x000000000021f030
  INIT_ARRAYSZ         0x0000000000000008
  FINI_ARRAY           0x000000000021f038
  FINI_ARRAYSZ         0x0000000000000008
  GNU_HASH             0x0000000000000298
  STRTAB               0x0000000000001080
  SYMTAB               0x0000000000000390
  STRSZ                0x00000000000005e8
  SYMENT               0x0000000000000018
  DEBUG                0x0000000000000000
  RELA                 0x00000000000017f0
  RELASZ               0x0000000000001db8
  RELAENT              0x0000000000000018
  BIND_NOW             0x0000000000000000
  FLAGS_1              0x0000000008000001
  VERNEED              0x0000000000001780
  VERNEEDNUM           0x0000000000000001
  VERSYM               0x0000000000001668
  RELACOUNT            0x00000000000000c2
:
```
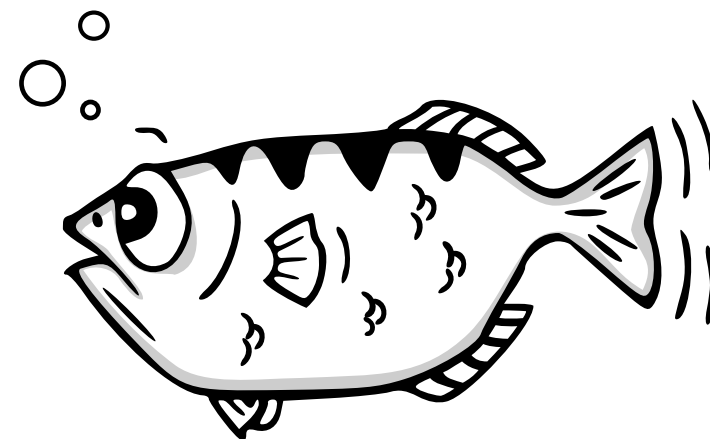
https://sourceware.org/binutils/

# Tracing

o Running the program and collecting information

  o Called Library Functions (with Arguments!) with 'ltrace'

  o Systemcalls (Files opened) with 'strace'

o Examples:

  o Binary deobfuscates some hostname and connects to it, so check for the 'connect' systemcalls

  o If some application just hangs the last syscall or library call might give you a hint

# Basic Runtime Influence: GDB

o GDB: GNU Debugger
  o Great for working with debugging symbols
  o Painful without them
o Can be enough for basic tasks on its own
  o Stop execution at certain addresses
  o Inspect registers and memory
o Plugins that help with analysis
  o https://github.com/longld/peda
  o https://github.com/pwndbg/pwndbg
  o https://github.com/hugsy/gef

Source: https://www.gnu.org/software/gdb/mascot/

https://www.gnu.org/software/gdb/

# Vanilla GDB vs Plugins

# Basic Runtime Influence

o LD_PRELOAD Functionality
  o Load your libraries before the specified ones
o Those Functions get called instead of the intended ones
  o Replace "getRandomNumber" with "rand"
  o gcc -shared -fPIC unrandom.c -o unrandom.so
  o LD_PRELOAD=$PWD/unrandom.so ./binary
o No more randomness!

```
int getRandomNumber()
{
    return 4;  // chosen by fair dice roll.
               // guaranteed to be random.
}
```

Source: https://xkcd.com/221/

# Intermediate

o Going deeper
  - o How do some tools work internally?
  - o Running non cooperative binaries in a controlled environment
  - o Specialized Tools

# Concept: Disassembly and Lifting

o  Map from bytes to an
   instruction

[0x83, 0xc0,0x01] -> "add eax,1"

# Concept: Disassembly and Lifting

o Map from bytes to an instruction

o Or lift to some other language that makes the semantics explicit

[0x83, 0xc0,0x01] -> "add eax,1"

add dstreg, immediate

dstreg += immediate

# Tool: Capstone

o Disassembly Framework

o Python (and other bindings)

o Many (FOSS) tools use it in the background somewhere


o Same project provides Keystone for assembly

https://www.capstone-engine.org/

# Disassembly Algorithms

o How do we or the tools know what part of the binary is code?

o ELF Information

   o Entrypoint
   o Possibly symbols

```
00000000: eb01 b848    ...H
00000004: c7c0 3905    ..9.
00000008: 0000 48c7    ..H.
0000000c: c37f 1d00    ....
00000010: 00ff d0      ...
```

# Disassembly Algorithms

o Linear Sweep
- o Easy to implement

o Might yield confusing results
- o On architectures like x86 with variable instruction lengths and no forced alignment

```
eb01:          jmp 3
b848c7c039:    mov eax, 0x39c0c748
05000048c7:    add eax, 0xc7480000
c3:            ret
7f1d:          jg 0x1f
0000:          add byte [rax], al
ffd0:          call rax
```

# Disassembly Algorithms

o **Recursive Descent Disassembly**
  - o Requires some semantic understanding
  - o More accurate

```
eb01:              jmp 3
b8:                db 0xb8
48c7c039050000:    mov rax, 1337
48c7c37f100000:    mov rbx, 4223
ffd0:              call rax
```

# Control Flow Graph Generation

o Graph of the possible control flows through the program

o Tradeoffs between accuracy and tractability

o Highly useful, it's easy to get lost in disassembly

o Every good graphical disassembler should have this somewhere

# Executable Parsing: Continued

- You might want to build something that needs this
  - Sure, you could just use objdump and grep

- Small pure python library for ELF parsing: 'pyelftools'

- If you want something more fancy: 'lief'

# Dynamic Analysis

o Observing and manipulating at runtime

| Target Process |
| --- |
| Environment(Filesystem/Libraries) |
| **Kernel** |
| **Execution Machine (CPU)** |

# Emulation

o Everything below some level of abstraction is emulated

o Level of Abstraction => Kind of emulation

| Target Process |
|:---:|

| Environment(Filesystem/Libraries) |
|:---:|

| Kernel |
|:---:|

| Execution Machine (CPU) |
|:---:|

# Emulation

| Target Binary | Target Binary | Target Binary |
|:---:|:---:|:---:|
| Libraries/Filesystem | Libraries/Filesystem | Libraries/Filesystem |
| Kernel | Kernel | Kernel |
| Machine | Machine | Machine |

Plain QEMU
User mode

Chrooting into
System Image with
QEMU as Interpreter

QEMU
System
Mode

# Emulation

- We are always at least in control of the execution machine
  - But we are slower than the real one
- Redefine instructions
  - Up is down, down is up, "inc reg" now decrements the register
- Add custom code to the emulation logic
  - Callback on every {jump, call, syscall} for analysis

- Fully emulating the environment might hard
  - Example: Windows API

# Tool: QEMU

- ○ Supports a lot of architectures
- ○ Used for device emulation in KVM/Xen
- ○ Decently fast
  - ○ JITs and caches basic blocks

Source: https://wiki.qemu.org/Logo

https://www.qemu.org/

# Tool: Unicorn

o Lightweight emulator
  o Just the CPU emulation core of QEMU
  o No device emulation
  o No syscalls
o Library
  o Use as the backend in some other tool
  o Emulate small code snippets

Source: https://www.unicorn-engine.org/images/unicorn.png

https://www.unicorn-engine.org/

# Tools: Misc

o **pyrebox**

   o IPython shell for introspection and instrumentation of (mainly Windows) guests

   o Main Focus: Malware Analysis

   o https://github.com/Cisco-Talos/pyrebox

o **panda2**

   o Full system tracing and analysis based on QEMU

   o https://github.com/panda-re/panda

# Dynamic Binary Instrumentation

o Rewrite the target code at runtime
  - o Remove code
  - o Add analysis code
  - o Hook functions

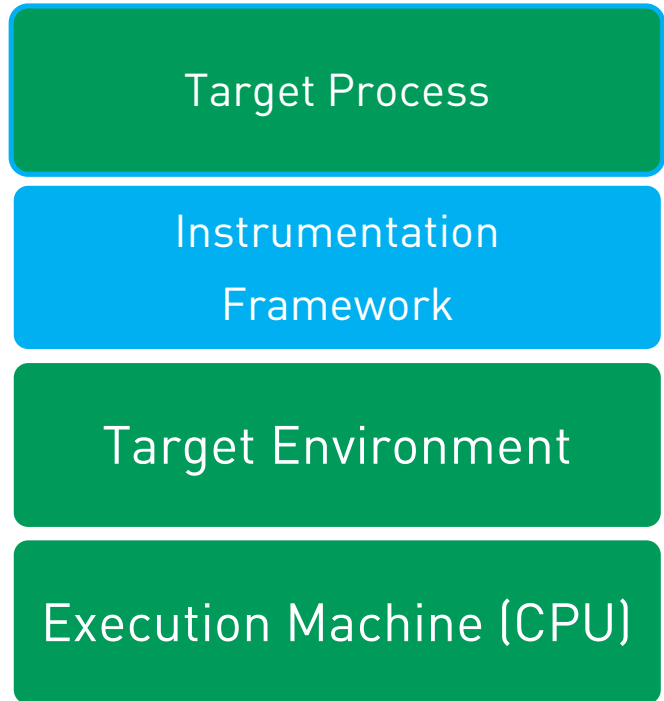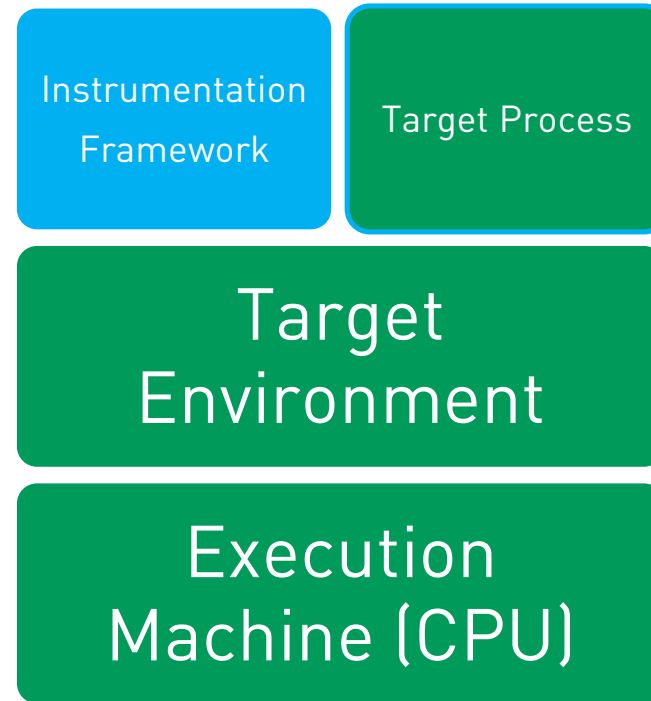| Target Process |
| :---: |
| Environment(Filesystem/Libraries) |
| **Kernel** |
| **Execution Machine (CPU)** |

# Some Dynamic Binary Instrumentation Approaches

| Target Process |
| :---: |

| Instrumentation Framework |
| :---: |

| Target Environment |
| :---: |

| Execution Machine (CPU) |
| :---: |

Framework runs a provided Binary.
Example: DynamoRIO

| Instrumentation Framework | Target Process |
| :---: | :---: |

| Target Environment | |
| :---: | :---: |

| Execution Machine (CPU) | |
| :---: | :---: |

Framework is loaded
into an existing process.
Example: Frida

# Dynamic Binary Instrumentation: Basic Idea



**BB1**
inc eax
jmp BB2

**BB2**
xor [edx], eax
jmp BB3

**BB3**
ret

Fig 1: Typical Execution Flow

**Instrumentation Core**
Per BB:
Parse next BB
Patch BB
Exec Block

**BB1**
inc eax
jmp core

**BB2**
xor [edx], eax
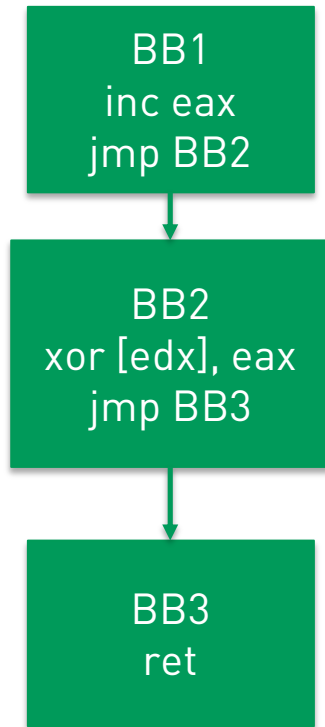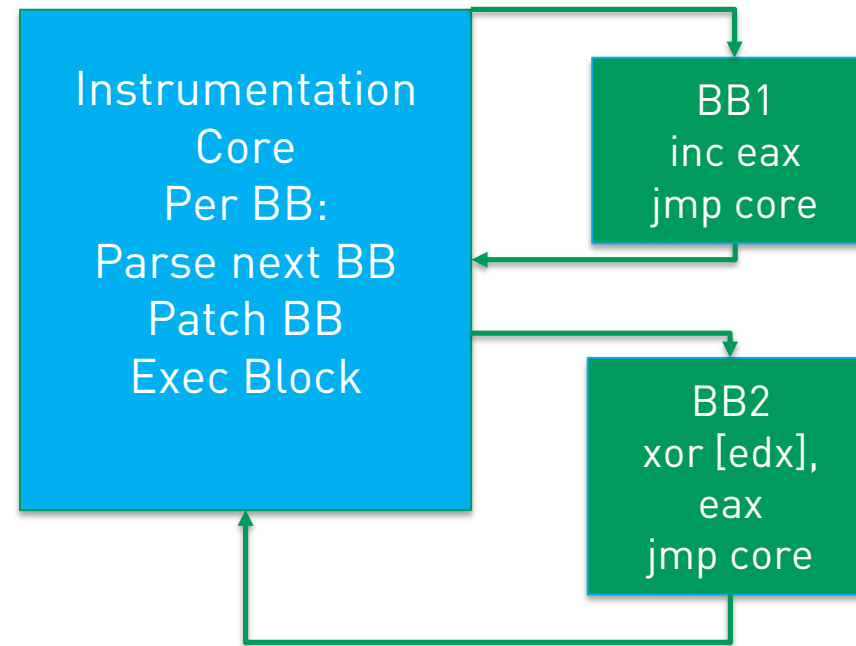jmp core

Fig 2: Execution Flow under Dynamic Instrumentation

# Dynamic Binary Instrumentation: Use Cases

o Hook functions
  o Library call and System Calling Tracing
  o Tracing of any function call
  o Basic Block Tracing for Coverage (Fuzzing)
  o Change return values

o Example
  o Static Crypto Key Generation is obfuscated?
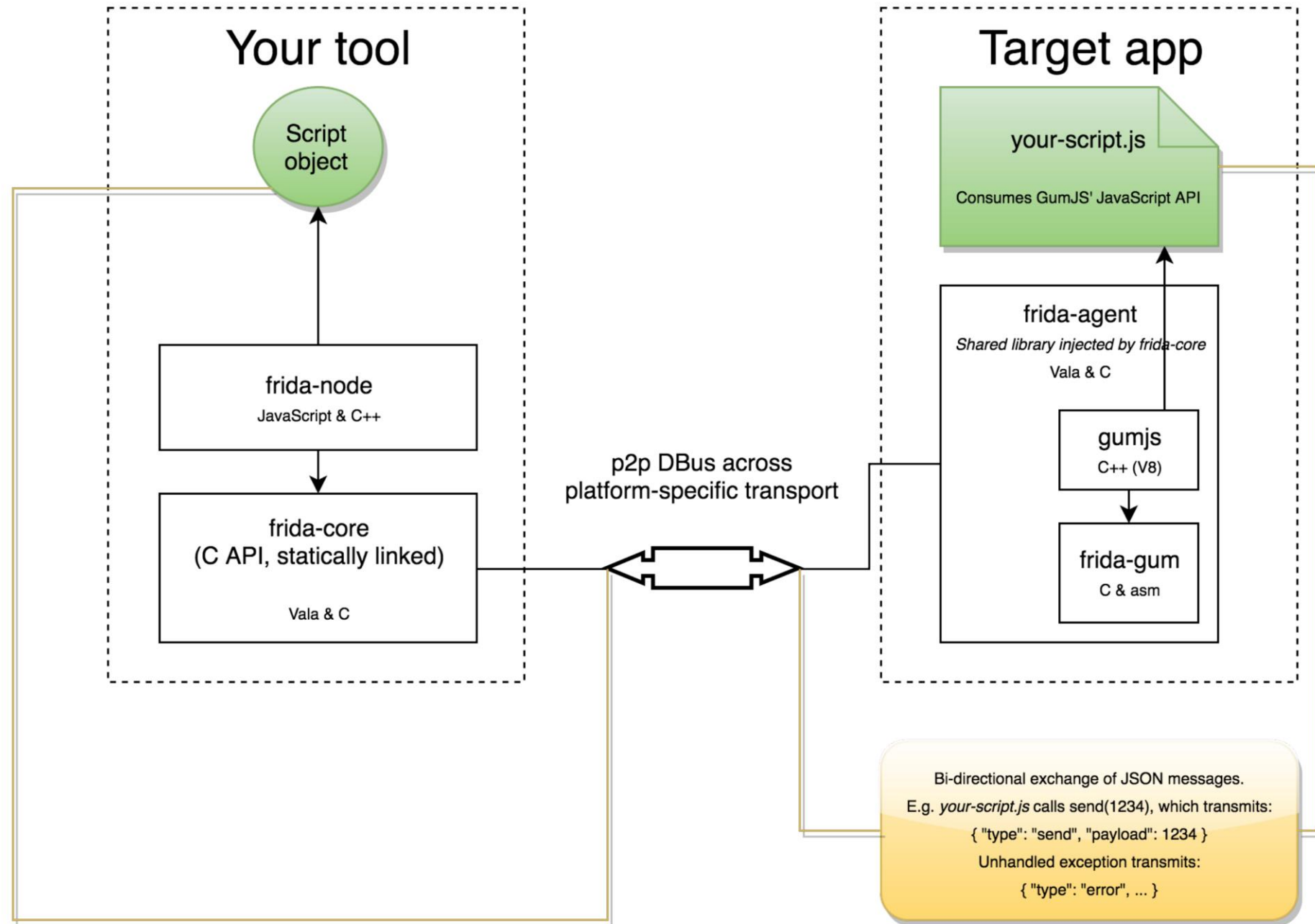  o Just hook the call where it is used

# Tool: Frida

FRIDA

Source: https://www.frida.re/img/logotype.svg

o Dynamic instrumentation toolkit

o Scriptable

o Multi-platform and multi-arch

   o Windows/Mac/Linux/Android/iOS/QNX – i386/AMD64/ARM/ARM64

o Bindings for Python, .NET, C and Node.js

   o But the actual scripts have to be written in Javascript...

o Very easy to Hook functions

https://www.frida.re/

Frida Architecture

# DBI: Misc

o DynamoRIO
  o More mature
  o FOSS (BSD)

o Intel PIN
  o More mature
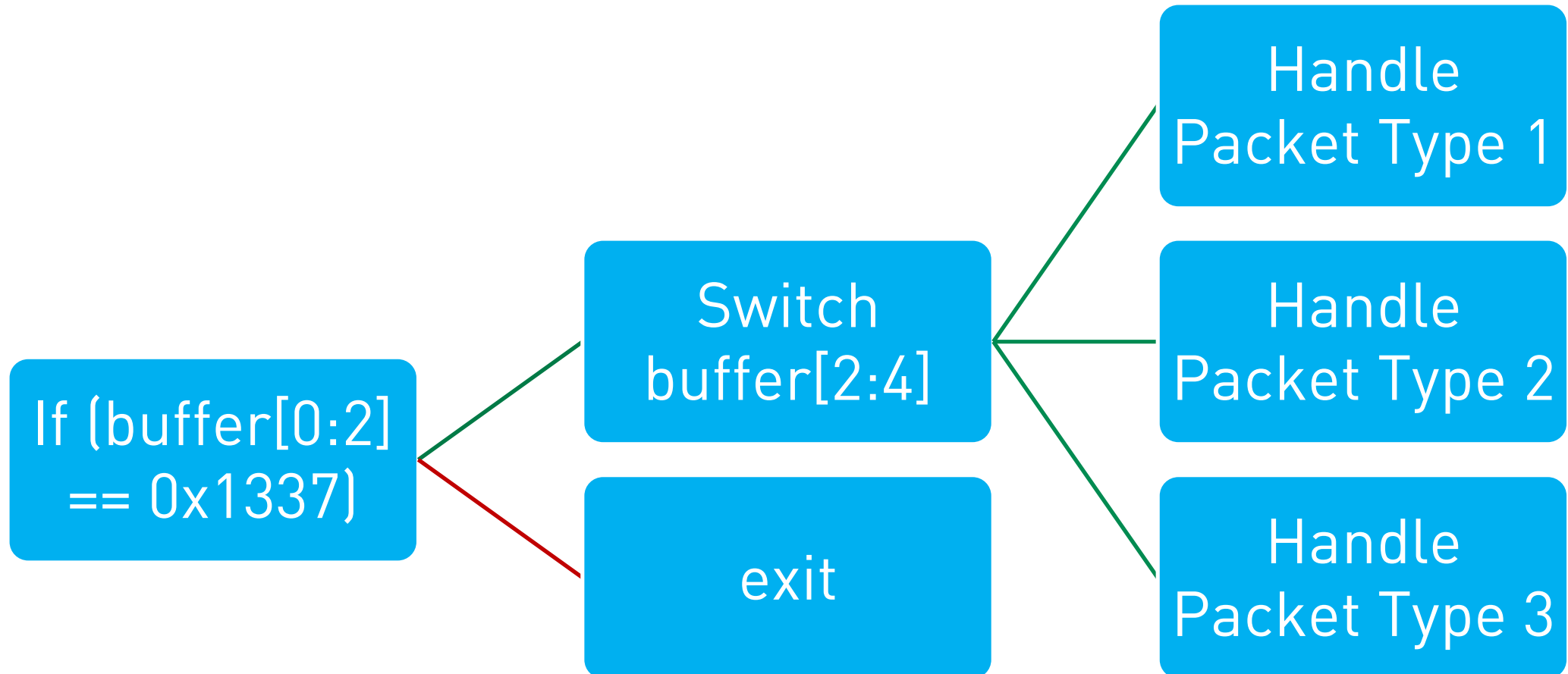  o Proprietary, but free as in beer

# Advanced

- ○ Formalizing
- ○ Automation

# Binary or Program Analysis

○ The subfield of computer science dealing with automated analysis

○ Massive improvements over the last years
  ○ Mainly due to the DARPA CGC

# Symbolic Execution



If (buffer[0:2] == 0x1337)

Switch buffer[2:4]

exit

Handle Packet Type 1

Handle Packet Type 2

Handle Packet Type 3

ERNW RESEARCH
pursuing knowledge.

# Tool: angr

o Binary Analysis Framework
  o Lifting to VEX IR
  o Emulation
  o Symbolic Execution
  o CFG Generation
o Used in the DARPA CGC by Shellphish , won 3rd place
o Best used from an interactive IPython Shell
o Build tools upon or integrate into others

Source: http://angr.io/img/angry_face.png

https://github.com/angr/angr

42

# Conclusion

o The right mature tooling makes your life a lot easier

o Initial learning overhead tends to be worth it


o Combine tooling to solve new problems


o Integrate new tooling into your existing tooling
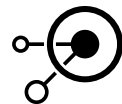
# Probably out of time?

Option 1: Q&A

Option 2: Misc

fmagin@ernw.de

@0x464D

www.ernw.de

www.insinuator.net

# Misc

# CPU Features

o CPUs sometimes provide advanced features

o Hardware Watchpoints
  o Hard to detect
  o Break on data read/write and not just code

o Intel PT
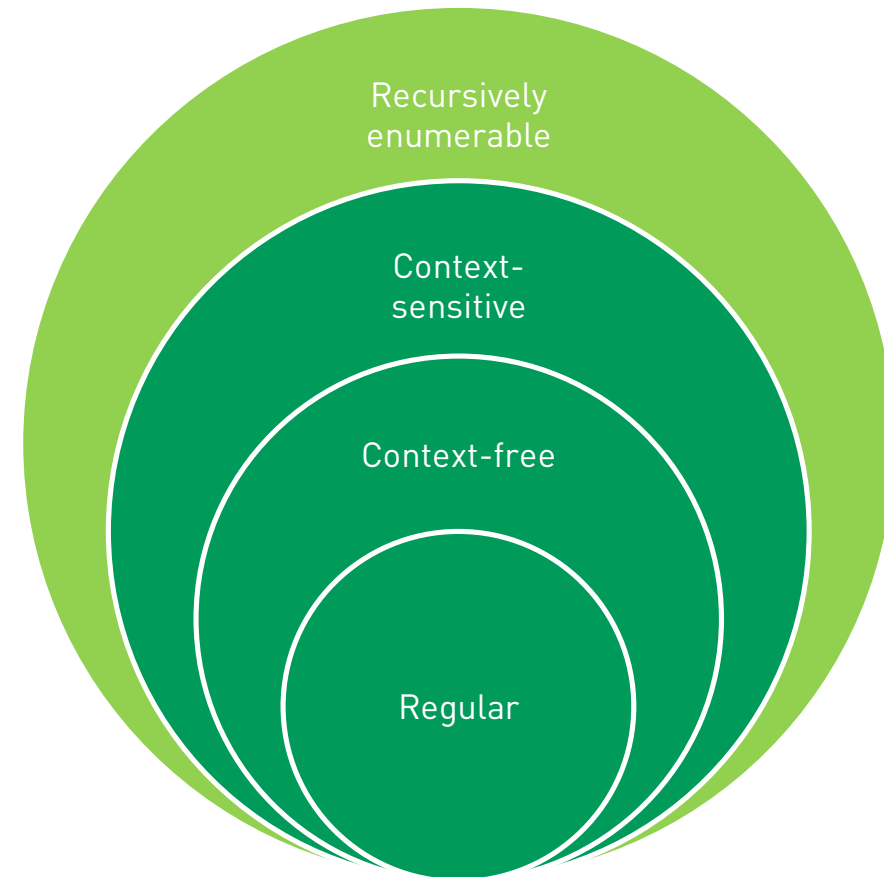  o Trace execution with your CPU

# Tool: rr

o Extension to GDB that allows recording a trace and debugging it

    o Run or step the program in reverse

https://rr-project.org/

# Concept: Parser Generation

o Unknown File Format

    o $proprietary protocol or file format

    o Some patterns might be obvious

    o Others can be derived from looking at an existing Parser

o Problem: Support for custom tools

    o Parser for Visualization

    o Parser/Serializer for Custom Client

    o Language Aware Fuzzer

# Theory: Formal Languages

- Every protocol or file format is basically a formal language
- Every formal language is induced by a grammar
  - More than one even
- You can generate a parser for the language from the grammar
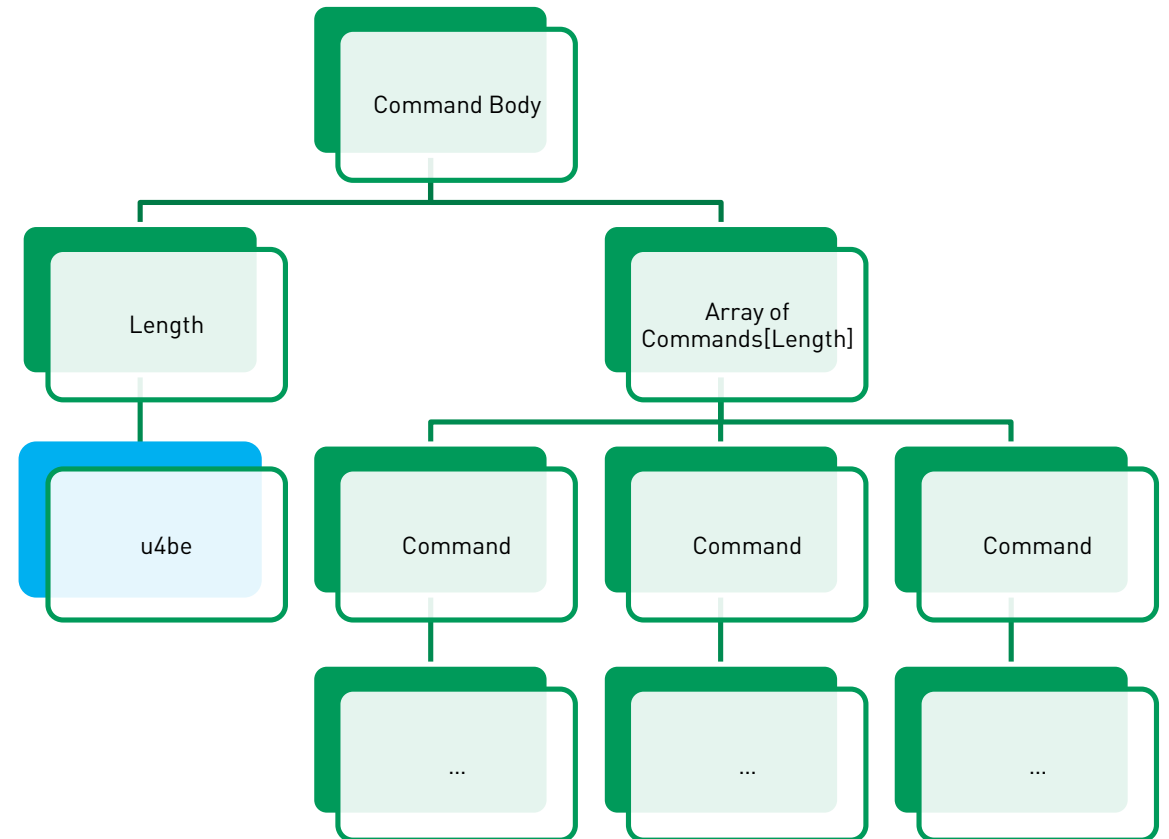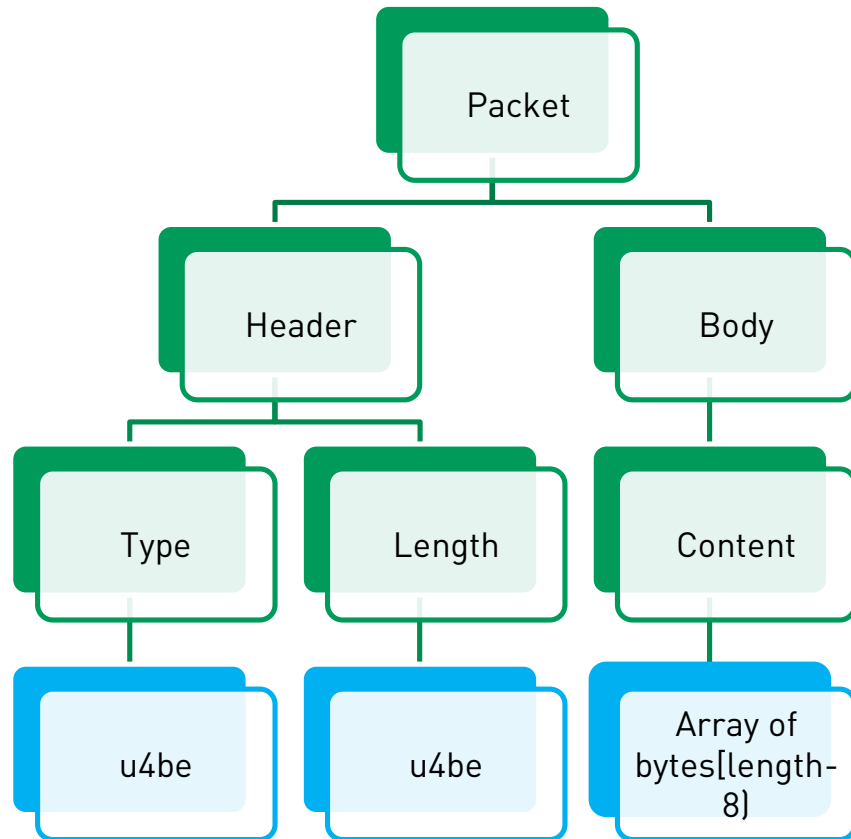  - Theoretically
- But for every sane protocol this should work

Recursively enumerable

Context-sensitive

Context-free

Regular

# Kaitai Struct

o Generates parser in many languages from a spec

o Compiles to
  o C++/STL
  o Python, Ruby, Perl
  o Javascript
  o C#, Java
  o Lua
  o Others

Source: http://kaitai.io/img/kaitai_16x_dark.png

# Parse Trees

# Kaitai Struct Use Cases

○ Binary protocols over HTTP that you are intercepting with 'mitmproxy'

○ Wireshark Dissectors

○ Burp Plugins