



Introduction to Binary Analysis and Binary Rewriting

Zhiqiang Lin

Department of Computer Science and Engineering
The Ohio State University

7/19/2018



Outline

- 1 Background
- 2 Challenges
- 3 Techniques
- 4 Applications
- 5 Summary

Outline

1 Background

2 Challenges

3 Techniques

4 Applications

5 Summary

What's Running in Our Cyberspace



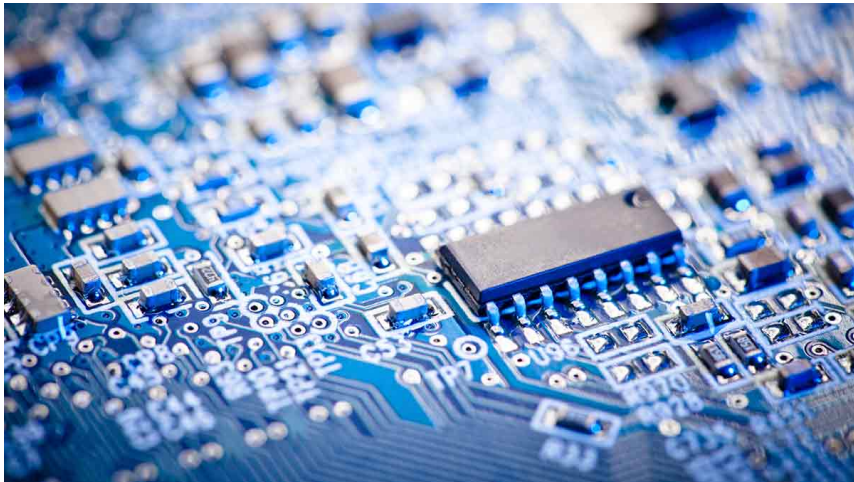
What's Running in Our Cyberspace

$$\begin{aligned}
 E_k &= \frac{1}{2} m v^2 \quad \tan \theta_B = \frac{v_2}{v_1} = v_{21} \quad \rho V = n R T \quad \psi = \iint \vec{D} d\vec{S} = A D \quad H_\lambda = \frac{\Delta E}{\Delta \lambda} \\
 -\frac{\hbar^2}{2m} \frac{d^2 \psi}{dx^2} + V \psi &= E \psi \quad M_e = \sigma T^4 \quad \phi_e = \frac{L}{\Delta t} \quad S = \frac{\Delta \varphi}{2\pi} = \frac{\Delta x}{\lambda} = \frac{x_2 - x_1}{\lambda} \quad v = c/\lambda \quad \Phi = NBS \\
 U_{ef} &= \frac{U_m}{E} \quad E = \hbar \omega \quad U = \frac{W_{AB}}{E_{PA} - E_{PB}} = \frac{|V_A - V_B|}{E_{PA} - E_{PB}} = \frac{|V_A - V_B|}{E_{PA} - E_{PB}} \quad X_L = \frac{U_m}{I_m} = \omega L = 2\pi f L \quad F_g = \frac{\mu_1 I_1 I_2}{2\pi d} \quad \vec{B} = \mu \frac{NI}{\sqrt{2}} \quad v = \frac{\omega \hbar}{m k} \quad \varphi_E = \frac{E_e}{E_0} = k \frac{\varphi}{r} \quad \varphi = \frac{Q}{4\pi \epsilon_0 r} \quad T = \frac{4 n_1 n_2}{(n_2 + n_1)^2} \quad g = \frac{\mu_1 \mu_2}{r^2} \quad \mathcal{R} = \frac{c}{T} k = \pm \sqrt{\frac{2m}{\hbar^2}} (E - V_0) \\
 K = \rho^2 / 2m \mu_0 = \frac{M_m}{N_A} = \frac{M_r \cdot 10^{-3}}{N_A} \quad l_t = l_0 (1 + d \Delta t) \quad I = \frac{U_e}{R + R_i} \quad E = \frac{E_c}{a} \int_0^a \sin(\omega t + \phi) dy \quad \omega = 2\pi f \\
 \lambda = \frac{h}{\sqrt{2e U m_e}} \quad R = \rho \frac{L}{S} \quad E = m c^2 \quad \frac{\sin \alpha}{\sin \beta} = \frac{v_1}{v_2} = \frac{\omega_2}{\omega_1} \quad v = \frac{1}{\sqrt{\epsilon \cdot \mu}} = \frac{c}{\sqrt{\epsilon_r \mu_r}} \\
 f_0 = \frac{1}{2\pi} \sqrt{\frac{g}{e}} \quad \psi_\alpha = \sqrt{2/L} \sin \frac{n\pi x}{L} \quad E = \frac{1}{2} \hbar / k_m \quad \beta = \frac{\Delta I_c}{\Delta t} \quad \phi_e = \frac{\Delta E}{\Delta t} \quad \frac{\omega_1}{x} + \frac{\omega_2}{x'} = \frac{\omega_2 - \omega_1}{r} \\
 \oint \vec{B} d\vec{l} = \mu \iint_S \vec{J} d\vec{S} \quad \vec{S} = \frac{1}{\mu_0} (\vec{E} \times \vec{B}) \quad E_k = \frac{h^2}{8mL^2} \quad \oint \vec{D} d\vec{S} = Q^* \\
 C(s) \quad v_k = \sqrt{\frac{3kT}{m_0}} = \sqrt{\frac{3kT N_A}{M_m}} = \sqrt{\frac{3R_m T}{M_r \cdot 10^{-3}}} \quad E = \hbar k^2 \quad 1 \text{ pc} = \frac{1 \text{ AU}}{r} \quad S = \frac{U}{I} \quad \psi_2 = U_e I t \quad F_v = \int \frac{F_n}{R} \\
 \lambda = \frac{h v_2}{T} \quad F_h = S h \rho g \quad f_0 = \frac{1}{2\pi \sqrt{LC}} \quad \sigma = \frac{Q}{M} \quad M = F d \cos \alpha \quad \lambda^* T = b \\
 \left(\frac{E_t}{E_0} \right)_\parallel = \frac{2 \cos \theta_1 \cos \theta_2}{\cos(\theta_1 - \theta_2) \sin(\theta_1 + \theta_2)} \quad \int \vec{E} d\vec{l} = - \iint_S \frac{\partial \vec{B}}{\partial t} \cdot d\vec{S} \quad p = \frac{E}{c} = \frac{h f}{c} = \frac{h}{\lambda} \quad \mu = U_m \sin \omega(t - T) = U_m \sin 2\pi \left(\frac{t}{T} - \frac{x}{\lambda} \right) \\
 S = \frac{1}{A} \frac{d\omega}{dt} \quad R = R_0 \sqrt[3]{A} \quad \int \vec{E} d\vec{l} = - \iint_S \frac{\partial \vec{B}}{\partial t} \cdot d\vec{S} \quad p = \frac{E}{c} = \frac{h f}{c} = \frac{h}{\lambda} \quad \mu = U_m \sin \omega(t - T) = U_m \sin 2\pi \left(\frac{t}{T} - \frac{x}{\lambda} \right)
 \end{aligned}$$

What's Running in Our Cyberspace



What's Running in Our Cyberspace



What's Running in Our Cyberspace



What is Binary Analysis

The process of (automatically) reasoning/deriving properties about the structure/behavior/syntactics/semantics/anything of your interest of binary programs

```
zlin@zlin-desktop:~/$ hexdump -C /bin/ls|less
00000000  7f 45 4c 46 02 01 01 00  00 00 00 00 00 00 00 00  |.ELF.....|
00000010  02 00 3e 00 01 00 00 00  a4 45 40 00 00 00 00 00  |..>.....E@...|
00000020  40 00 00 00 00 00 00 00  70 96 01 00 00 00 00 00  |@.....p.....|
00000030  00 00 00 00 40 00 38 00  09 00 40 00 1c 00 1b 00  |....@.8...@...|
00000040  06 00 00 00 05 00 00 00  40 00 00 00 00 00 00 00  |.....@.....|
00000050  40 00 40 00 00 00 00 00  40 00 40 00 00 00 00 00  |@.@.....@.@...|
00000060  f8 01 00 00 00 00 00 00  f8 01 00 00 00 00 00 00  |.....|
00000070  08 00 00 00 00 00 00 00  03 00 00 00 04 00 00 00  |.....|
00000080  38 02 00 00 00 00 00 00  38 02 40 00 00 00 00 00  |8.....8.@....|
00000090  38 02 40 00 00 00 00 00  1c 00 00 00 00 00 00 00  |8.@.....|
...
```

Why Binary Code?

Access to the source code often is not possible:

- ▶ Proprietary software packages. (Volks Wagon's cheating software)
- ▶ Stripped executables.
- ▶ Proprietary libraries
- ▶ Malicious software (exploits), e.g., stuxnet

Why Binary Code?

Access to the source code often is not possible:

- ▶ Proprietary software packages. (Volks Wagon's cheating software)
- ▶ Stripped executables.
- ▶ Proprietary libraries
- ▶ Malicious software (exploits), e.g., stuxnet

Binary code is the only authoritative version of the program.

- ▶ Binary code is everywhere
- ▶ Changes occurring in the compile, optimize and link steps can create non-trivial semantic differences from the source and binary.
- ▶ What you see is not what you execute (WYSINWYX problem)

Why Binary Code?


- ▶ Windows
 - ▶ Login process keeps a user's password in the heap after a successful login
- ▶ To minimize data lifetime
 - ▶ clear buffer
 - ▶ call `free()`

```
memset(buffer, '\\0', len);  
free(buffer);
```

Why Binary Code?

- ▶ Windows
 - ▶ Login process keeps a user's password in the heap after a successful login
- ▶ To minimize data lifetime
 - ▶ clear buffer
 - ▶ call free()
- ▶ But . . .
 - ▶ the compiler might optimize away the buffer-clearing code ("useless-code" elimination)

```
memset(buffer, '\0', len);  
free(buffer);
```



Why Binary Code: Backdoor



Linux embedded device: HTTP server for management and video monitoring, with a known backdoor.

Backdoor!!!

- Username: 3sadmin
- Password: 27988303

```
LDR    R1, =a3sadmin ; "3sadmin"
MOV    R0, R7 ; s1
BL     strcmp
CMP    R0, #0
LDR    R1, =a27988303 ; "27988303"
MOV    R0, R4 ; s1
```

Heffner, Craig. "Finding and Reversing Backdoors in Consumer Firmware." EELive! (2014).

What to Reason About in Binary Code?

The process of (automatically) reasoning/deriving properties about the structure/behavior/syntactics/semantics/anything of your interest of binary programs

- ❶ What are the program's variables and their types?
- ❷ What are the program's parameters and their types?
- ❸ Where could this indirect jump go?
- ❹ What function could be called at this indirect call site?
- ❺ What could this dereference operation access/affect?
- ❻ What kind of object is allocated at this allocation site?
- ❼ What could the value held in V eventually affect?
- ❽ What could have affected the value of V ?
- ❾ What are the statements (at instruction level) that contribute to the execution of i ?
- ❿ ...

Outline

- 1 Background
- 2 Challenges
- 3 Techniques
- 4 Applications
- 5 Summary

Challenges: Abstraction Recovery

The first step in any binary code analysis is to **reconstruct the abstractions distilled after compilation**, such as recognizing the instructions, operand, opcode, variables, basic blocks, and control flows

Challenges: Abstraction Recovery

The first step in any binary code analysis is to **reconstruct the abstractions distilled after compilation**, such as recognizing the instructions, operand, opcode, variables, basic blocks, and control flows

Challenges

- ▶ Code/Data distinction
- ▶ Variable x86 instruction size
- ▶ Indirect Branches
- ▶ Functions without explicit CALL
- ▶ Position independent code (PIC)
- ▶ ...

It will be easier to recover these abstractions by using **dynamic analysis**, but will be much more challenge in **static analysis**.

Challenges: Path Coverage, and Path Explosion

For both static analysis and dynamic analysis, how to model the program execution path (too conservative, or too simple), and how to trigger the program path (especially for dynamic analysis) is another key challenge

Challenges: Path Coverage, and Path Explosion

For both static analysis and dynamic analysis, how to model the program execution path (too conservative, or too simple), and how to trigger the program path (especially for dynamic analysis) is another key challenge

Static analysis

- ▶ Too conservative
- ▶ Too many paths
- ▶ Impossible path

Dynamic analysis

- ▶ A single path
- ▶ Cover more path
- ▶ Test case generation

A Surprise:

Many source-level issues gone

- ① Use of multiple source languages
- ② In-lined assembly code
- ③ Avoid building problems
- ④ Analyze the actual libraries
- ⑤ ...

Many people inspecting binaries

- ① IDA Pro Users
- ② Anti-malware companies
- ③ Computer Emergency Response Teams
- ④ Malware writers
- ⑤ ...

A Surprise:

- ▶ Many exploits utilize platform-specific quirks
 - ▶ non-obvious and unexpected
 - ▶ compiler artifacts (choices made by the compiler)
 - ▶ Memory layout
 - padding between fields of a struct
 - which variables are adjacent
 - ▶ register usage
 - ▶ execution order
 - ▶ optimizations performed
 - ▶ compiler bugs

Outline

- 1 Background
- 2 Challenges
- 3 Techniques**
- 4 Applications
- 5 Summary

Basic Approaches to disassemble code

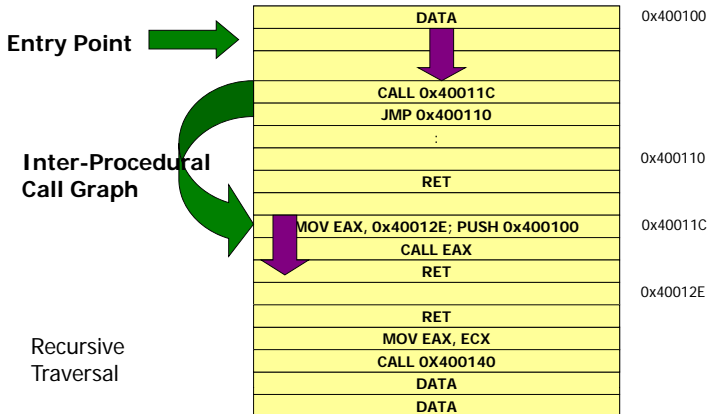
Linear sweep algorithm

- ▶ Start with program entry point, proceed to disassemble instructions sequentially
- ▶ Key assumption: all instructions appear one after the next, without any gaps
 - ▶ Violated in most code (presence of data or padding)

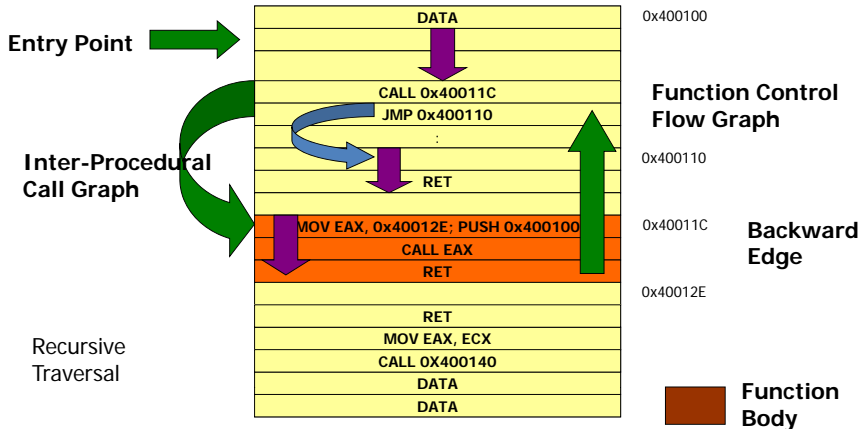
Recursive Traversal Algorithm

- ▶ After a control-flow transfer instruction (CTI), proceed to disassemble target address
- ▶ For conditional CTI and non-CTI, proceed to disassemble next instruction
- ▶ Key problems
 - ▶ Code reached only through indirect CTIs
 - ▶ Functions that do not return in the usual way

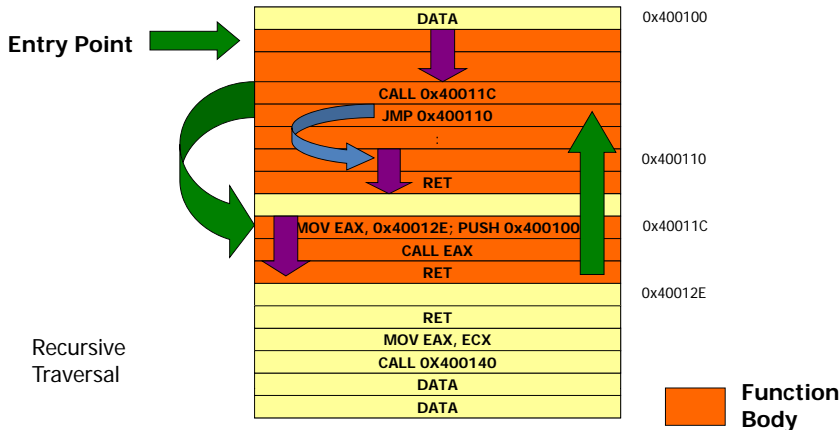
Examples



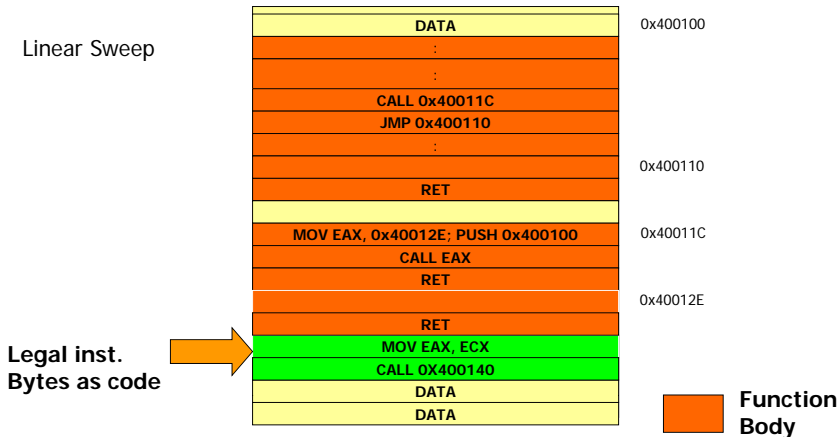
Examples



Examples



Examples



Examples

Error Recovery Heuristics

Code sequence
not ending with
JMP/RET
→ not code



DATA	0x400100
:	
:	
CALL 0x40011C	
JMP 0x400110	
:	
	0x400110
RET	
MOV EAX, 0x40012E; PUSH 0x400100	0x40011C
CALL EAX	
RET	
	0x40012E
RET	
MOV EAX, ECX	
CALL 0x400140	
DATA	
DATA	

 **Function
Body**

Disassemble package

- ❶ Udis86 is an easy-to-use, minimalistic disassembler library (libudis86) for the x86 class of instruction set architectures <http://udis86.sourceforge.net/>.
- ❷ libdasm, a disassembly library. <https://code.google.com/p/libdasm/>.
- ❸ pydasm <https://github.com/OpenRCE/pydbg>
- ❹ Capstone is a lightweight multi-platform, multi-architecture disassembly framework. <http://www.capstone-engine.org/>
- ❺ ...

Basic Techniques

- ❶ Data Flow Analysis
 - ▶ Data dependency
 - ▶ Taint analysis
 - ▶ Point-to analysis
- ❷ Control Flow Analysis
 - ▶ Control flow graph
 - ▶ Call graph
 - ▶ Control dependency
- ❸ Program Slicing
- ❹ Symbolic Execution

Static Analysis, Dynamic Analysis, Symbolic Execution

BAP BAT CodeReason
radare2
vivisect Hex-Ray IDA rdis Valgrind
amoco fuzzgrind gdb
angr SemTrax
BitBlaze
BARF JARVIS
klee/s2e Jakstab
insight
PIN QEMU Bindead
Triton
PySysEmu TEMU PEMU
miasm CodeSurfer
paimei

Common Tools

- ❶ Static analysis
 - ▶ IDA Pro, BinNav
 - ▶ BAP
- ❷ Dynamic analysis
 - ▶ PIN
 - ▶ QEMU
 - ▶ PEMU
- ❸ Symbolic execution
 - ▶ FuzzBall, Fuzzgrind
 - ▶ S2E
 - ▶ Angr

Outline

- 1 Background
- 2 Challenges
- 3 Techniques
- 4 Applications**
- 5 Summary

Applications of Binary Analysis in Security

Use Cases

- ➊ Reverse engineering (knowing the secrets)
- ➋ Vulnerability discovery/fuzzing
- ➌ Exploit generation
- ➍ Software verification
- ➎ Program testing
- ➏ Binary hardening
- ➐ ...

Applications: Vulnerability Discovery

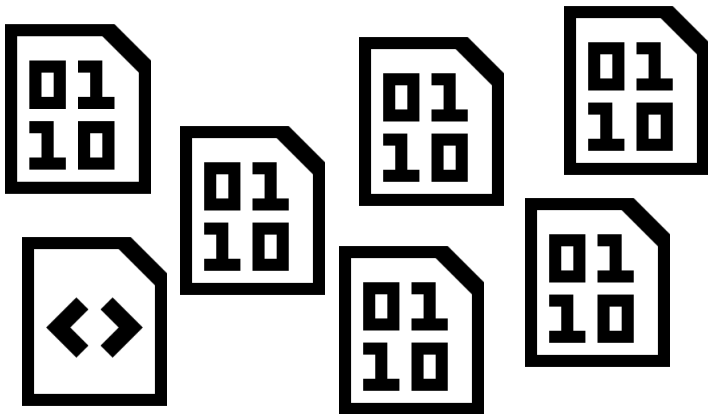


Vulnerability Discovery

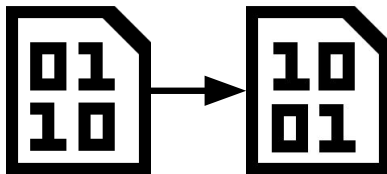
Basic Approaches

- ❶ Static Analysis
 - ▶ “You can’t ” / “You might be able to”
 - ▶ Based on various static techniques.
- ❷ Dynamic Analysis
 - ▶ Input A? Input B? Input C? ...
 - ▶ Based on concrete inputs to application
- ❸ Symbolic Execution
 - ▶ Interpret the application.
 - ▶ Track “constraints” on variables.
 - ▶ When the required condition is triggered, “concretize” to obtain a possible input

Applications: Binary Rewriting



Applications: Binary Rewriting



Applications: Binary Rewriting

- ➊ Software fault isolation (SFI)
- ➋ Control Flow Integrity (CFI)
- ➌ Binary code hardening
- ➍ Binary code reuse
- ➎ Platform-specific optimizations
- ➏ Vulnerability fuzzing by rewriting binary code

Outline

- 1 Background
- 2 Challenges
- 3 Techniques
- 4 Applications
- 5 Summary

Binary Analysis

- ▶ Binary code is everywhere, and it is the final representation of programs
- ▶ Binary analysis is challenging
- ▶ It is extremely useful to perform binary code analysis (vulnerability excavation, backdoor identification) in security
- ▶ Basic binary analysis approaches: static/dynamic analysis, symbolic execution
- ▶ There are many public available binary analysis tools