

Xenpwn

Breaking Paravirtualized Devices

Felix Wilhelm





#whoami



- Security Researcher @ ERNW Research
- Application and Virtualization Security
- Recent Research
 - Security Appliances (Fireeye, Palo Alto)
 - Hypervisors (Xen)
- ¬ @_fel1x on Twitter



Agenda

- Device Virtualization & Paravirtualized Devices
- Double Fetch Vulnerabilities
- Xenpwn: Architecture and Design
- Results
- Case Study: Exploiting xen-pciback



Device Virtualization

intel، 82078 44 PIN CHMOS SINGLE-CHIP FLOPPY DISK CONTROLLER

- Small Footprint and Low Height Package
- Enhanced Power Management
- Application Software Transparency
 Programmable Powerdown
- Command — Save and Restore Commands for
- Zero-Volt Powerdown — Auto Powerdown and Wakeup
- Modes — Two External Power Management
- Pins
- Consumes No Power While in Powerdown
- Integrated Analog Data Separator
- 250 Kbps
- 300 Kbps
- 500 Kbps
- 1 Mbps
- Programmable Internal Oscillator
- Floppy Drive Support Features
- Drive Specification Command
- Selectable Boot Drive
- Standard IBM and ISO Format Features
- Format with Write Command for High Performance in Mass Floppy Duplication

- Integrated Tape Drive Support
 Standard 1 Mbps/500 Kbps/ 250 Kbps Tape Drives
- Perpendicular Recording Support for 4 MB Drives
- Integrated Host/Disk Interface Drivers
- Fully Decoded Drive Select and Motor Signals
- Programmable Write Precompensation Delays
- Addresses 256 Tracks Directly, Supports Unlimited Tracks
- 16 Byte FIFO
- Single-Chip Floppy Disk Controller Solution for Portables and Desktops 100% DC (ATt Competition)
- 100% PC/AT* Compatible
- Fully Compatible with Intel386™ SL
 Integrated Drive and Data Bus Buffers
- Separate 5.0V and 3.3V Versions of the 44 Pin part are Available
- Available in a 44 Pin QFP Package

- Virtualized systems need access to virtual devices
 - Disk, Network, Serial, ...
- Traditionally: Device emulation
 - Emulate old and well supported hardware devices
 - Guest OS does not need special drivers
 - Installation with standard installation sources supported



Paravirtualized Devices

- Most important downsides of emulated devices:

- Hard to implement securely and correctly
- Slow performance
- No support for advanced features

- Solution: Paravirtualized Devices

- Specialized device drivers for use in virtualized systems
- Idea: Emulated devices are only used as fallback mechanism
- Used by all major hypervisors



Paravirtualized Devices



- Split Driver Model

- Frontend runs in Guest system
- Backend in Host/Management domain
- Terminology differs between hypervisors
 - VSC/VSP in Hyper-V
 - Virtio devices and drivers
- Implementations are quite similar



Paravirtualized Devices



PV devices are implemented on top of shared memory

- Great Performance
- Easy to implement
- Zero copy algorithms possible
- Message protocols implemented on top
 - Xen, Hyper-V and KVM all use ring buffers
- Shared memory mappings can be constant or created on demand



Security of PV Devices

- ¬ Backend runs in privileged context → Communication between frontend and backend is trust boundary
- ¬ Low level code + Protocol parsing → Bugs
- Examples
 - Heap based buffer overflow in KVM disk backend (CVE-2011-1750)
 - Unspecified BO in Hyper-V storage backend (CVE-2015-2361)
- Not as scrutinized as emulated devices
 - Device and hypervisor specific protocols
 - Harder to fuzz



Very interesting target

- ¬ Device emulation often done in user space ← → PV backend often in kernel for higher performance
 - Compromise of kernel backend is instant win 🕲
- PV devices are becoming more important
 - More device types (USB, PCI pass-through, touch screens, 3D acceleration)
 - More features, optimizations
- Future development: Removal of emulated devices
 - see Hyper-V Gen2 VMs



Research goal

- "Efficient vulnerability discovery in Paravirtualized Devices"
- Core Idea: No published research on the use of shared memory in the context of PV devices
- ¬ Bug class that only affect shared memory? → Double fetches!



Double Fetch vulnerabilities

- ¬ Special type of TOCTTOU bug affecting shared memory.
- Simple definition: Same memory address is accessed multiple times with validation of the accessed data missing on at least one access
- Can introduce all kinds of vulnerabilities
 - Arbitrary Write/Read
 - Buffer overflows
 - Direct RIP control☺



Double Fetch vulnerabilities

6	1	2	2
134	15	20	2
12	13	34	2
	100	1	6

Felix Wilhelm @_fel1x - Dec 17 Blogpost about some interesting double fetch vulnerabilities I discovered in Xen: insinuator.net/2015/12/xen-xs... #XSA155

★ 17 65 ♥ 52 ···

grsecurity arsecurity

👤 Follow

@_fel1x Small history fix: twitter.com/grsecurity/sta..., see also: osronline.com/article.cfm?ar...

grsecurity @grsecurity "double fetch" in 2007 by sgrakkyu/twiz (see 2.4.2); phrack.org/issues.html?is...

- Term "double fetch" was coined by Fermin J. Serna in 2008
 - But bug class was well known before that
- Some interesting research published in 2007/2008
 - Usenix 2007 "Exploiting Concurrency Vulnerabilities in System Call Wrappers" -Robert N. M. Watson
 - CCC 2007: "From RING 0 to UID 0" and Phrack #64 file 6 – twiz, sgrakkyu
- First example I could find is sendmsg() linux bug reported in 2005
 - Happy to hear about more ☺

```
1 int cmsghdr_from_user_compat_to_kern(..)
2 {
    [...]
3
    while(ucmsg != NULL) {
4
          if(get_user(ucmlen, &ucmsg->cmsg_len))
                   return -EFAULT;
6
          [...]
7
          tmp = ((ucmlen - CMSG_COMPAT_ALIGN(sizeof(*ucmsg))) +
8
                 CMSG_ALIGN(sizeof(struct cmsghdr)));
9
          kcmlen += tmp;
10
          [...]
11
    }
12
13
    if(kcmlen > stackbuf_size)
14
          kcmsg_base = kcmsg = kmalloc(kcmlen, GFP_KERNEL);
15
16
    while(ucmsg != NULL) {
17
                                                                      Example: sendmsg()
          __get_user(ucmlen, &ucmsg->cmsg_len);
19
          if(copy_from_user(CMSG_DATA(kcmsg),
20
                   CMSG_COMPAT_DATA(ucmsg),
21
           (ucmlen - CMSG_COMPAT_ALIGN(sizeof(*ucmsg)))))
22
23 [...]
24 }
4/14/16
```





Bochspwn



- "Identifying and Exploiting Windows Kernel Race Conditions via Memory Access Patterns" (2013)
 - by j00ru and Gynvael Coldwind
- Uses extended version of Bochs
 CPU emulator to trace all memory access from kernel to user space.



Bochspwn

- Resulted in significant number of Windows bugs (and a Pwnie)
 - but not much published follow-up research
- Whitepaper contains detailed analysis on exploitability of double fetches
 - On multi core system even extremely short races are exploitable
- Main inspiration for this research.



1	mov	ecx, [edi+18h]	
2	;[]		
3	push	4	
4	push	eax	
5	push	ecx	
6	call	_ProbeForWrite	
7	push	dword ptr [esi+20h]	
8	push	dword ptr [esi+24h]	Example: Bochspwn
9	push	dword ptr [edi+18h]	ntlApphalpCachal aaku
10	call	_memcpy	pEntry



Xenpwn

 Adapt memory access tracing approach used by Bochspwn for analyzing PV device communication.

- Why not simply use Bochspwn?

- Extremely slow
- Passive overhead (no targeted tracing)
- Compatibility issues
- Dumping traces to text files does not scale
- Idea: Implement memory access tracing on top of hardware assisted virtualization







Xenpwn Architecture



¬ Nested virtualization

- Target hypervisor (L1) runs on top of base hypervisor (L0)
- Analysis components run in user space of L1 management domain.
 - No modification to hypervisor required
 - Bugs in these components do not crash whole system
- LO hypervisor is Xen



LibVMI



- Great library for virtual machine introspection (VMI)
 - Hypervisor agnostic (Xen and KVM)
 - User-space wrapper around hypervisor APIs
- Allows access to and manipulation of guest state (memory, CPU registers)
- Xen version supports memory events



LibVMI Memory Events

```
auto event = new vmi_event_t();
event->type = VMI_EVENT_MEMORY;
event->mem_event.physical_address = paddr;
event->mem_event.npages = 1;
event->mem_event.granularity = granularity;
event->mem_event.in_access = access;
event->callback = callback;
```

```
if (vmi_register_event(s->vmi, event) != VMI_SUCCESS)
{ /*... */}
```

- Trap on access to a guest physical address
- Implemented on top of Extended Page Tables (EPT)
 - Disallow access to GPA
 - Access triggers EPT violation and VM exit
 - VM exit is forwarded to libvmi handler



Memory Access Tracing with libVMI



- 1. Find shared memory pages
- 2. Register memory event handlers
- 3. Analyze memory event, extract needed information and store in trace storage.
- 4. Run analysis algorithms (can happen much later)



Trace Collector

- Use libvmi to inspect memory and identify shared memory pages
 - Target specific code.
 - Identify data structures used by PV frontend/backend and addresses of shared pages
- Registers memory event handlers
- Main work is done in callback handler
 - Disassemble instructions using Capstone







Trace Storage

- Storage needs to be fast and persistent

- Minimize tracing overhead
- Allow for offline analysis
- Nice to have: Efficient compression
 - Allows for very long traces
- Tool that fulfills all these requirements: Simutrace
 - simutrace.org



Simutrace

Simutrace

- Open source project by the Operation System Group at the Karlsruhe Institute of Technology
- Designed for full system memory tracing
 - All memory accesses including their content
- C++ daemon + client library
 - Highly efficient communication over shared memory pages
- Uses specialized compression algorithm optimized for memory traces
 - High compression rate + high speed
- Highly recommended!



Trace Entries

For every memory access:



For every unique instruction:







Simplified version (Ignores overlapping accesses and interweaved read/writes)



Advantages & Limitations

- Good:

- Low passive overhead
- Largely target independent
 - only Trace collector requires adaption
- Easy to extend and develop
- Bad
 - High active overhead
 - VM exits are expensive
 - Reliance on nested virtualization



Nested Virtualization on Xen

- Xen Doku: Nested HVM on Intel CPUs, as of Xen 4.4, is considered "tech preview". For many common cases, it should work reliably and with low overhead
- Reality:
 - Xen on Xen works
 - KVM on Xen works (most of the time)
 - Hyper-V on Xen does not work $\ensuremath{\mathfrak{S}}$
- For this reason, all of the following results are from Xen
 - .. but still hopeful for Server 2016 Hyper-V



Results

Tracing runs for two L1 targets:

Component	Xen-Ubuntu	Xen-SLES
L1 Hypervisor	Xen 4.5.0	Xen 4.4.2_08-1.7
L2 dom0 OS	Ubuntu 15.04	SLES 11 SP4
L2 dom0 Kernel	3.19.0-18-generic	3.0.101-63-xen
Management Stack	xl	xend

- Differences in supported PV devices
 - SCSI, USB



Results

- Main Problem: Getting good coverage

- No automated way to exercise device functionality implemented
- In the following: Interesting bugs found with default compiler settings
 - Full thesis contains more statistic about instruction types and attack surface

```
void blkif_get_x86_64_req(blkif_request_t *dst,
                           blkif_x86_64_request_t *src)
        int i, n = BLKIF MAX SEGMENTS PER REQUEST;
        dst->operation = src->operation;
        dst->nr_segments = src->nr_segments;
        // ...
        if (src->operation == BLKIF_OP_DISCARD) {
                //..
        }
        if (n > src->nr_segments)
                                                      QEMU xen disk
                 n = src->nr_segments;
        for (i = 0; i < n; i++)</pre>
                                                      Normally not exploitable thanks
                 dst->seg[i] = src->seg[i];
                                                      to compiler optimizations
```





```
1 \text{ for } (n = 0, i = 0; n < nseg; n++) {
       //...
2
        i = n % SEGS PER INDIRECT FRAME;
3
        seg[n].nsec = segments[i].last_sect -
                 segments[i].first_sect + 1;
5
6
        seg[n].offset = (segments[i].first_sect << 9);</pre>
7
8
        if ((segments[i].last_sect >= (PAGE_SIZE >> 9)) ||
9
            (segments[i].last_sect < segments[i].first_sect)) {</pre>
10
                rc = -EINVAL;
11
                                                             xen-blkback
                goto unmap;
12
                                                             00B Read/Write
13
       11...
14
15 }
```



xen-pciback







xen-pciback: xen_pcibk_do_op

```
1 switch (op->cmd) {
       case XEN_PCI_OP_conf_read:
2
                op->err = xen pcibk config read(dev,
3
                          op->offset, op->size, &op->value);
                break;
5
       case XEN_PCI_OP_conf_write:
6
               11...
7
       case XEN_PCI_OP_enable_msi:
8
               //...
9
       case XEN PCI OP disable msi:
10
              //...
11
       case XEN_PCI_OP_enable_msix:
12
                   11...
13
       case XEN_PCI_OP_disable_msix:
14
              11...
15
       default:
16
                op->err = XEN PCI ERR not implemented;
17
                break;
18
19 }
```

1 Cmp	DWORD PTR [r13+0x4],0x5
2 mov	DWORD PTR [rbp-0x4c],eax
з ја	0x3358 <xen_pcibk_do_op+952></xen_pcibk_do_op+952>
4 MOV	eax, <mark>DWORD</mark> PTR [r13+0x4]
5 jmp	QWORD PTR [rax*8+off_77D0]



xen-pciback

1	cmp	DWORD PTR [r13+0x4],0x5
2	mov	DWORD PTR [rbp-0x4c],eax
3	ja	0x3358 <xen_pcibk_do_op+9523< th=""></xen_pcibk_do_op+9523<>
4	mov	eax, DWORD PTR [r13+0x4]
5	jmp	QWORD PTR [rax*8+off_77D0]

- switch statement is compiled into jump table
- ¬ op->cmd == \$r13+0x4
- Points into shared memory
- Range check and jump use two different memory accesses
- Valid compiler optimization
 - op is not marked as volatile



Exploiting pciback

1	cmp	DWORD PTR [r13+0x4],0x5
2	mov	DWORD PTR [rbp-0x4c],eax
3	ja	0x3358 <xen_pcibk_do_op+952></xen_pcibk_do_op+952>
4	mov	eax,DWORD PTR [r13+0x4]
5	jmp	QWORD PTR [rax*8+off_77D0]

```
"loop_header_%=:\n"
"inc rcx\n"
"xor dword ptr [rax], 25\n"
"cmp rcx, 5000\n"
"jnz loop_header_%=\n"
```

- Race is very small: 2 Instructions

- But can be reliably won if guest VM has multiple cores
- Lost race does not have any negative side effects
 - Infinite retries possible

- Simple to trigger

 Send PCI requests while flipping value using XOR



Exploiting pciback

- ¬ Indirect jump → No immediate RIP control
 - Need to find reliable offset to function pointer
- Load address of xen-pciback.ko is random
- Virtual address of backend mapping also not known
- A lot of similarities to a remote kernel exploit
- Chosen approach: Trigger type confusion to get write primitive



Type Confusion

<pre>xen_pcibk_frontend_changed(struct xenbus_device *xdev,</pre>
enum xenbus_state fe_state)
<pre>struct xen_pcibk_device *pdev = dev_get_drvdata(&xdev->dev);</pre>
<pre>switch (fe_state) {</pre>
case XenbusStateInitialised:
<pre>xen_pcibk_attach(pdev);</pre>
break;
case XenbusStateReconfiguring:
<pre>xen_pcibk_reconfigure(pdev);</pre>
break;
//
//

- 1 mov rdi, r13
- 2 call 0x3720 <xen_pcibk_attach>

Second jump table generated for xen-pciback

- Almost directly behind the jump table generated for vulnerable function
- XenbusStateInitialized uses value of r13 register as first argument
 - Should be a pointer to a xen_pcibk_device structure
 - Is a pointer to the start of the shared memory page ☺



Getting a write primitive

struct xen_pcibk_device {
 void *pci_dev_data;
 struct mutex dev_lock;
 struct xenbus_device *xdev;
 struct xenbus_watch be_watch;
 u8 be_watching;
 int evtchn_irq;
 struct xen_pci_sharedinfo *sh_info;
 unsigned long flags;
 struct work_struct op_work;
 struct xen_pci_op op;

};

void __sched mutex_lock(struct mutex *lock)

```
might_sleep();
```

/*

* The locking fastpath is the 1->0 transition from

```
* 'unlocked' into 'locked' state.
```

```
__mutex_fastpath_lock(&lock->count, __mutex_lock_slowpath);
mutex_set_owner(lock);
```

- xen_pcibk_attach first tries to lock the dev_lock mutex of referenced structure.
- Gives us the possibility to call mutex_lock with a fake mutex structure
- mutex_lock
 - Fastpath: Switch lock count from 1 -> 0
 - Slowpath: Triggered when lock count != 1



Getting a write primitive: mutex_lock slowpath

1.

/* add waiting tasks to the end of the waitqueue (FIF0): */
list_add_tail(&waiter.list, &lock->wait_list);
waiter.task = task;

```
wait_list->prev = new;
waiter->next = wait_list;
waiter->prev = WRITE_TARGET;
WRITE TARGET->next = new;
```

mutex_optimistic_spin needs to fail.

- Can be achieved by setting lock->owner to a readable zero page
- If lock count still not 1, mutex_waiter
 structure is created and stored on stack
- mutex_waiter structure is added to lock->wait_list and kernel thread goes to sleep till wake up.
- Pointer to waiter is written to attacker controlled location.



Write Primitive

struct list_head {
 struct list_head *next, *prev;
};

write-where but not write-what

- Pointer to pointer to attacker controlled data
- Can't simply overwrite function pointers

- One shot

- pciback is locked due to xen_pcibk_do_op never returning
- Idea: Add faked entries to a global linked list.
 - Requires known kernel version + no KASLR or infoleak





4/14/16

#44 www.ernw.de







Overwrite Target



- Global data structure

- Need to know address of list_head
- No new elements should be attached during run time
 - list_head.prev is not changed, new entry might be added directly behind list_head
- Needs to survive one "junk" entry
 - No full control over waiter structure / stack frame



```
/*
1
2
3
4
5
6
7
    *
       linux/fs/exec.c
    *
       Copyright (C) 1991, 1992 Linus Torvalds
    */
   /*
8
9
     #!-checking implemented by tytso.
    */
10
   /*
11
    * Demand-loading implemented 01.12.91 - no need to read anything but
12
    * the header into memory. The inode of the executable is put into
13
14
    * "current->executable", and page faults do the actual loading. Clean.
15
    * Once more I can proudly say that linux stood up to being changed: it
16
17
    * was less than 2 hours work to get demand-loading completely implemented.
    *
18
    * Demand loading changed July 1993 by Eric Youngdale. Use mmap instead,
19
    * current->executable is only used by the procfs. This allows a dispatch
20
    * table to check for several different types of binary formats. We keep
21
    * trying until we recognize the file or we run out of supported binary
22
23
    * formats.
    */
```



fs/exec.c: formats

static LIST_HEAD(formats);

```
list_for_each_entry(fmt, &formats, lh) {
    if (!try_module_get(fmt->module))
        continue;
    read_unlock(&binfmt_lock);
    bprm->recursion_depth++;
    retval = fmt->load_binary(bprm);
    read_lock(&binfmt_lock); -
```

- formats linked list contains entries for different file formats supported by exec
 - ELF
 - #! shell scripts
 - a.out format
- Walked every time exec* syscall is called to load input file.
- waiter entry is skipped because try_module_get function fails



Getting Code Execution

- Set address of load_binary pointer to stack pivot
- ROP chain to allocate executable memory and copy shellcode
 - vmalloc_exec + memcpy
- Restore original formats list
- \$shellcode
- Return to user space



Demo 🙂







Thesis, Whitepaper & Code

- Master Thesis describing Xenpwn in greater detail can be found online:
 - https://os.itec.kit.edu/downloads/ma_2015_wilhelm_felix ______discover__software__vulnerabilities.pdf
- Exploit code + Whitepaper for pciback vulnerability will be released after Infiltrate
- Xenpwn open source release: May 2016



Future Work

- Use Xenpwn against Hyper-V and VMWare

- Requires improved support for nested virtualization
- Identify and analyze other shared memory trust boundaries
 - Sandboxes?
- What types of bugs can we find with full memory traces?



Thanks for your attention!







fwilhelm@ernw.de

ld fel1x



