# Advanced Exploitation:
# Xen Hypervisor VM Escape

Shangcong Luan

May 27, 2016

Alibaba Cloud Platform Security Team

2014 - 2015   at Vulnhunt Security Team for APT Defense

2015 - now   at Alibaba Cloud Platform Security Team for Cloud Security
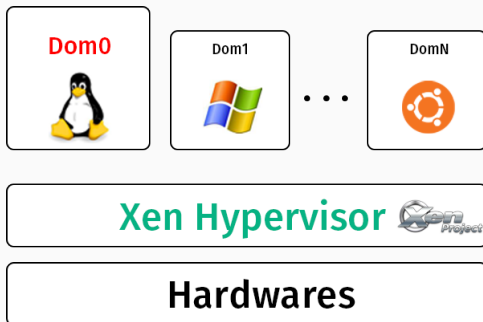
Twitter:   @hikerell

1. Introduction

2. XSA-148/CVE-2015-7835

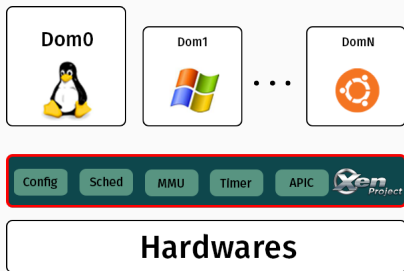3. Exploitation Technologies

4. The End

5. Demo Time

# Introduction

*" The Xen Project$^{TM}$ is the leading open source virtualization platform that is powering some of the largest clouds in production today. "*

*from xenproject.org*

- CPU Scheduling
- Memory Management
- VM Executation
- ...

Dom0:

- Privileged Domain
- Control Other Domains

domU:

- Dom1, Dom2, Dom3 ...
- Unprivileged Domains

PVM:

- paravirtualization machine
- modified OS kernel

HVM:

- hardware-assisted virtualization machine
- unmodified OS kernel
- CPU/MMU => hardware assistance

x86 Paravirtualised Memory Management:

Direct Paging

machine frame number (mfn)

||

guest pseudo-physical frame number (gpfn)

mutually-exclusive page types:

- PGT_writable_page could be writable mapped into Guest
- PGT_l1_page_table L1 page table type
- PGT_l2_page_table L2 page table type
- PGT_l3_page_table L3 page table type
- PGT_l4_page_table L4 page table type

1. A guest OS may always create readable mappings to its own page frames, regardless of their current types.

2. A frame may only safely be retasked when its reference count is zero.

PV Guest Cannot Read/Write Security-Sensitive Memorys, e.g., page tables

1. A guest OS may always create readable mappings to its own page frames, regardless of their current types.

2. A frame may only safely be retasked when its reference count is zero.

PV Guest Cannot Read/Write Security-Sensitive Memorys, e.g., page tables

1.  A guest OS may always create readable mappings to its own page frames, regardless of their current types.

2.  A frame may only safely be retasked when its reference count is zero.

PV Guest Cannot Read/Write Security-Sensitive Memorys,

e.g., page tables.

XSA-148/CVE-2015-7835

## Offical Vulnerablity Advisoriy[1]

### Information

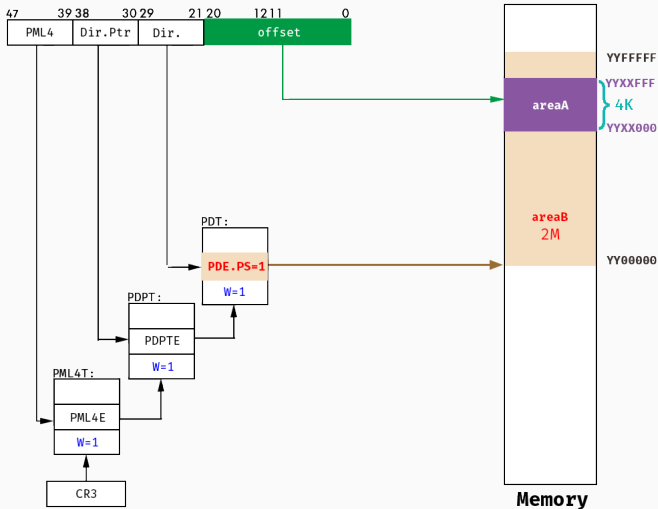| | |
|---|---|
| Advisory | XSA-148 |
| Public release | 2015-10-29 11:59 |
| Updated | 2015-10-29 11:59 |
| Version | 4 |
| CVE(s) | CVE-2015-7835 |
| Title | x86: Uncontrolled creation of large page mappings by PV guests |

### Files

advisory-148.txt (signed advisory file)
xsa148.patch
xsa148-4.4.patch
xsa148-4.5.patch

---

[1]http://xenbits.xen.org/xsa/advisory-148.html

```
xen-4.4.0 xen/arch/x86/mm.c
1756 /* Update the L2 entry at pl2e to new value nl2e. pl2e is within frame pfn. */
1757 static int mod_l2_entry(l2_pgentry_t *pl2e,
1758                         l2_pgentry_t nl2e,
1759                         unsigned long pfn,
1760                         int preserve_ad,
1761                         struct vcpu *vcpu)
1762 {
1763     l2_pgentry_t ol2e;
1764     struct domain *d = vcpu->domain;
1765     struct page_info *l2pg = mfn_to_page(pfn);
1766     unsigned long type = l2pg->u.inuse.type_info;
1767     int rc = 0;
1768
        check-1:
1769     if ( unlikely(!is_guest_l2_slot(d, type, pgentry_ptr_to_slot(pl2e))) )
1770     {
1771         MEM_LOG("Illegal L2 update attempt in Xen-private area %p", pl2e);
1772         return -EPERM;
1773     }
1774
1775     if ( unlikely(__copy_from_user(&ol2e, pl2e, sizeof(ol2e)) != 0) )
1776         return -EFAULT;
1777
```

# Page Tables Update Validation

```
         check-2: PDE.P == 1 ?
1778     if ( l2e_get_flags(nl2e) & _PAGE_PRESENT )
1779     {
             check-3: PDE.reserved_bits == 0 ?
1780         if ( unlikely(l2e_get_flags(nl2e) & L2_DISALLOW_MASK) )
1781         {
1782             MEM_LOG("Bad L2 flags %x",
1783                     l2e_get_flags(nl2e) & L2_DISALLOW_MASK);
1784             return -EINVAL;
1785         }
1786
             check-4: Old PDE.PADDR == New PDE.PADDR and Old PDE.P == New PDE.P ?
1787         /* Fast path for identical mapping and presence. */
1788         if ( !l2e_has_changed(ol2e, nl2e, _PAGE_PRESENT) )
1789         {
1790             adjust_guest_l2e(nl2e, d);
1791             if ( UPDATE_ENTRY(l2, pl2e, ol2e, nl2e, pfn, vcpu, preserve_ad) )
1792                 return 0;
1793             return -EBUSY;
1794         }
1795
         // check-5: do other audit, such as super page audit
1796         if ( unlikely((rc = get_page_from_l2e(nl2e, pfn, d)) < 0) )
1797             return rc;
         // ...
1806     }
     // ...
1815 }
```

# Page Tables Update Validation

```
          check-2: PDE.P == 1 ?
1778      if ( l2e_get_flags(nl2e) & _PAGE_PRESENT )
1779      {
              check-3: PDE.reserved_bits == 0 ?
1780          if ( unlikely(l2e_get_flags(nl2e) & L2_DISALLOW_MASK) )
1781          {
1782              MEM_LOG("Bad L2 flags %x",
1783                      l2e_get_flags(nl2e) & L2_DISALLOW_MASK);
1784              return -EINVAL;
1785          }
1786

              check-4: Old PDE.PADDR == New PDE.PADDR and Old PDE.P == New PDE.P ?
1787          /* Fast path for identical mapping and presence. */
1788          if ( !l2e_has_changed(ol2e, nl2e, _PAGE_PRESENT) )
1789          {
1790              adjust_guest_l2e(nl2e, d);
1791              if ( UPDATE_ENTRY(l2, pl2e, ol2e, nl2e, pfn, vcpu, preserve_ad) )
1792                  return 0;
1793              return -EBUSY;
1794          }
1795

          // check-5: do other audit, such as super page audit
1796          if ( unlikely((rc = get_page_from_l2e(nl2e, pfn, d)) < 0) )
1797              return rc;
          // ...
1806      }
      // ...
1815  }
```

```
        check-2: PDE.P == 1 ?
1778    if ( l2e_get_flags(nl2e) & _PAGE_PRESENT )
1779    {
            check-3: PDE.reserved_bits == 0 ?
1780        if ( unlikely(l2e_get_flags(nl2e) & L2_DISALLOW_MASK) )
1781        {
1782            MEM_LOG("Bad L2 flags %x",
1783                    l2e_get_flags(nl2e) & L2_DISALLOW_MASK);
1784            return -EINVAL;
1785        }
1786
            check-4: Old PDE.PADDR == New PDE.PADDR and Old PDE.P == New PDE.P ?
1787        /* Fast path for identical mapping and presence. */
1788        if ( !l2e_has_changed(ol2e, nl2e, _PAGE_PRESENT) )
1789        {
1790            adjust_guest_l2e(nl2e, d);
1791            if ( UPDATE_ENTRY(l2, pl2e, ol2e, nl2e, pfn, vcpu, preserve_ad) )
1792                return 0;
1793            return -EBUSY;
1794        }
1795
        // check-5: do other audit, such as super page audit
1796        if ( unlikely((rc = get_page_from_l2e(nl2e, pfn, d)) < 0) )
1797            return rc;
        // ...
1806    }
    // ...
1815 }
```
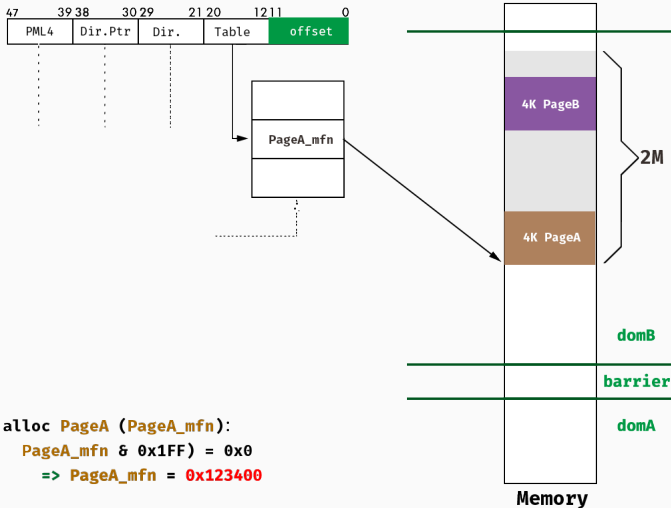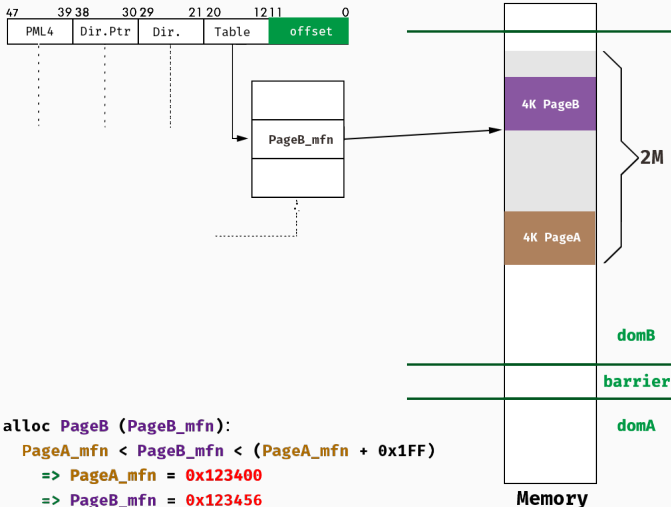
How to trigger it?

translate the XSA-148 to

Arbitrary Physical Memory Read/Write

```
alloc PageA (PageA_mfn):
  PageA_mfn & 0x1FF) = 0x0
    => PageA_mfn = 0x123400
```

```
alloc PageB (PageB_mfn):
    PageA_mfn < PageB_mfn < (PageA_mfn + 0x1FF)
        => PageA_mfn = 0x123400
        => PageB_mfn = 0x123456
```

PageB filled with 0
PageB registered as a page table

PageA filled with 0
PageA registered as a page table

# Exploitation Technologies

Hijack Vectors:

- Hypercall Page
- VDSO/vsyscall Page
- Hypercall Table
- ...

Hijack Vectors:

- Hypercall Page
- VDSO/vsyscall Page
- Hypercall Table
- ...

- A 4K page alloced by Guest Kernel filled with 0xCC
- Need to be initialized with hypercall stub codes
- Guest Kernel use it to do hypercall request
- Each Guest Kernel only has one hypercall page

- A 4K page alloced by Guest Kernel filled with 0xCC
- Need to be initialized with hypercall stub codes
- Guest Kernel use it to do hypercall request
- Each Guest Kernel only has one hypercall page

- A 4K page alloced by Guest Kernel filled with 0xCC
- Need to be initialized with hypercall stub codes
- Guest Kernel use it to do hypercall request
- Each Guest Kernel only has one hypercall page

- A 4K page alloced by Guest Kernel filled with 0xCC
- Need to be initialized with hypercall stub codes
- Guest Kernel use it to do hypercall request
- Each Guest Kernel only has one hypercall page

```
569   static void hypercall_page_initialise_ring3_kernel(void *hypercall_page)
570 ▼ {
571       char *p;
572       int i;
573       /* Fill in all the transfer points with template machine code. */
574       for ( i = 0; i < (PAGE_SIZE / 32); i++ )
575 ▼     {
576           if ( i == __HYPERVISOR_iret )
577               continue;
578           p = (char *)(hypercall_page + (i * 32));
579           *(u8  *)(p+ 0) = 0x51;      /* push %rcx */
580           *(u16 *)(p+ 1) = 0x5341;    /* push %r11 */
581           *(u8  *)(p+ 3) = 0xb8;      /* mov  $<i>,%eax */
582           *(u32 *)(p+ 4) = i;
583           *(u16 *)(p+ 8) = 0x050f;    /* syscall */
584           *(u16 *)(p+10) = 0x5b41;    /* pop  %r11 */
585           *(u8  *)(p+12) = 0x59;      /* pop  %rcx */
586           *(u8  *)(p+13) = 0xc3;      /* ret */
587       }
588 ▼     /*
589        * HYPERVISOR_iret is special because it doesn't return and expects a
590        * special stack frame. Guests jump at this transfer point instead of
591        * calling it.
592        */
593       p = (char *)(hypercall_page + (__HYPERVISOR_iret * 32));
594       *(u8  *)(p+ 0) = 0x51;      /* push %rcx */
595       *(u16 *)(p+ 1) = 0x5341;    /* push %r11 */
596       *(u8  *)(p+ 3) = 0x50;      /* push %rax */
597       *(u8  *)(p+ 4) = 0xb8;      /* mov  $__HYPERVISOR_iret,%eax */
598       *(u32 *)(p+ 5) = __HYPERVISOR_iret;
599       *(u16 *)(p+ 9) = 0x050f;    /* syscall */
600   }
```

36

```
569  static void hypercall_page_initialise_ring3_kernel(void *hypercall_page)
570 ▼ {
571      char *p;
572      int i;
573      /* Fill in all the transfer points with template machine code. */
574      for ( i = 0; i < (PAGE_SIZE / 32); i++ )
575 ▼   {
576          if ( i == __HYPERVISOR_iret )
577              continue;
578          p = (char *)(hypercall_page + (i * 32));
579          *(u8  *)(p+ 0) = 0x51;      /* push %rcx */
580          *(u16 *)(p+ 1) = 0x5341;    /* push %r11 */
581          *(u8  *)(p+ 3) = 0xb8;      /* mov  $<i>,%eax */
582          *(u32 *)(p+ 4) = i;
583          *(u16 *)(p+ 8) = 0x050f;    /* syscall */
584          *(u16 *)(p+10) = 0x5b41;    /* pop  %r11 */
585          *(u8  *)(p+12) = 0x59;      /* pop  %rcx */
586          *(u8  *)(p+13) = 0xc3;      /* ret */
587      }
588 ▼   /*
589       * HYPERVISOR_iret is special because it doesn't return and expects a
590       * special stack frame. Guests jump at this transfer point instead of
591
592
593
594
595
596
597
598
599
600  }
```

**Hypercall Page Signature:**

0x00000000B8534151

0xCCCCC3595B41050F

37

+0x0000        <- hypercall page stub codes begin

... 

+0x02E0        _iret stub

syscall

...

+0x0800        <- hypercall page stub codes end

...        <- unused 2048 bytes

+0x0FFF

38

+0x0000     <- hypercall page stub codes begin

...

+0x02E0
            _iret stub
                            <- replaced with 'JMP xxx'
            syscall

...

            <- hypercall page stub codes end
+0x0800
            <- hook_restore codes

            <- ret2usermode codes

            <- fork codes

            <- cxt restore codes

            <- shell codes

            <- BUFFER
+0x0FFF

```
569  static void hypercall_page_initialise_ring3_kernel(void *hypercall_page)
570  {
571      char *p;
572      int i;
573      /* Fill in all the transfer points with template machine code. */
574      for ( i = 0; i < (PAGE_SIZE / 32); i++ )
575      {
576          if ( i == __HYPERVISOR_iret )
577              continue;
578          p = (char *)(hypercall_page + (i * 32));
579          *(u8  *)(p+ 0) = 0x51;    /* push %rcx */
580          *(u16 *)(p+ 1) = 0x5341;  /* push %r11 */
581          *(u8  *)(p+ 3) = 0xb8;    /* mov  $<i>,%eax */
582          *(u32 *)(p+ 4) = i;
583          *(u16 *)(p+ 8) = 0x050f;  /* syscall */
584          *(u16 *)(p+10) = 0x5b41;  /* pop  %r11 */
585          *(u8  *)(p+12) = 0x59;    /* pop  %rcx */
586          *(u8  *)(p+13) = 0xc3;    /* ret */
587      }
588      /*
589       * HYPERVISOR_iret is special because it doesn't return and expects a
590       * special stack frame. Guests jump at this transfer point instead of
591       * calling it.
592       */
593      p = (char *)(hypercall_page + (__HYPERVISOR_iret * 32));
594      *(u8  *)(p+ 0) = 0x51;    /* push %rcx */
595      *(u16 *)(p+ 1) = 0x5341;  /* push %r11 */
596      *(u8  *)(p+ 3) = 0x50;    /* push %rax */
597      *(u8  *)(p+ 4) = 0xb8;    /* mov  $_HYPERVISOR_iret,%eax */
598      *(u32 *)(p+ 5) = __HYPERVISOR_iret;
599      *(u16 *)(p+ 9) = 0x050f;  /* syscall */
600  }
```

## Hypercall Page:

```
+00*32:    pushq %rcx
           pushq %r11
           movq $0x0, %rax
           syscall
           popq %r11
           popq %rcx
           ...
+23*32:    push %rcx
           push %r11
           push %rax
           mov $_HYPERVISOR_iret, %rax
           syscall
           ...
```
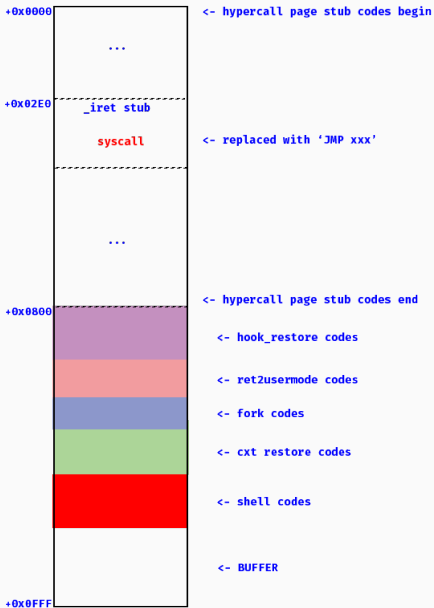
```
569  static void hypercall_page_initialise_ring3_kernel(void *hypercall_page)
570 ▼ {
571      char *p;
572      int i;
573      /* Fill in all the transfer points with template machine code. */
574      for ( i = 0; i < (PAGE_SIZE / 32); i++ )
575 ▼   {
576          if ( i == __HYPERVISOR_iret )
577              continue;
578          p = (char *)(hypercall_page + (i * 32));
579          *(u8  *)(p+ 0) = 0x51;    /* push %rcx */
580          *(u16 *)(p+ 1) = 0x5341;  /* push %r11 */
581          *(u8  *)(p+ 3) = 0xb8;    /* mov  $<i>,%eax */
582          *(u32 *)(p+ 4) = i;
583          *(u16 *)(p+ 8) = 0x050f;  /* syscall */
584          *(u16 *)(p+10) = 0x5b41;  /* pop  %r11 */
585          *(u8  *)(p+12) = 0x59;    /* pop  %rcx */
586          *(u8  *)(p+13) = 0xc3;    /* ret */
587      }
588 ▼   /*
589       * HYPERVISOR_iret is special because it doesn't return and expects a
590       * special stack frame. Guests jump at this transfer point instead of
591       * calling it.
592       */
593      p = (char *)(hypercall_page + (__HYPERVISOR_iret * 32));
594      *(u8  *)(p+ 0) = 0x51;    /* push %rcx */
595      *(u16 *)(p+ 1) = 0x5341;  /* push %r11 */
596      *(u8  *)(p+ 3) = 0x50;    /* push %rax */
597      *(u8  *)(p+ 4) = 0xb8;    /* mov  $_HYPERVISOR_iret,%eax */
598      *(u32 *)(p+ 5) = __HYPERVISOR_iret;
599      *(u16 *)(p+ 9) = 0x050f;  /* syscall */
600  }
```

## Hypercall Page:

+00*32:
```
pushq %rcx
pushq %r11
movq $0x0, %rax
syscall
popq %r11
popq %rcx
...
```

+23*32:
```
push %rcx
push %r11
push %rax
mov $_HYPERVISOR_iret, %rax
syscall JMP XXXX
...
```

41

# Hypercall Page Hijacked Layout



+0x0000    <- hypercall page stub codes begin

...

+0x02E0    _iret stub

           syscall          <- replaced with 'JMP xxx'

...

+0x0800    <- hypercall page stub codes end

           <- hook_restore codes

           <- ret2usermode codes

           <- fork codes

           <- cxt restore codes

           <- shell codes

           <- BUFFER

+0x0FFF

42

# The End

$ git show f87f8a7110e5dd57091b8484685953414693e2a3

```
Date:   Tue Feb 8 15:13:45 2005 +0000

+
+    if ( l2_pgentry_val(nl2e) & _PAGE_PRESENT )
+    {
+        /* Differ in mapping (bits 12-31) or presence (bit 0)? */
+        if ( ((l2_pgentry_val(ol2e) ^ l2_pgentry_val(nl2e)) & ~0xffe) == 0 )
+            return update_l2e(pl2e, ol2e, nl2e);
+
```

- dom0, dom1, dom2 ...
- PV or HVM guest ...
- linux or windows ...

# Demo Time

**Host:** Debian 7 with Xen 4.4.0

**Guest:** Debian 8

**Video:** Get Dom0 Shell (XSA-148)

Question?

_____

Shangcong Luan from Alibaba Cloud Platform Security Team