

Basic Concepts of Virtualization

Junming Liu

Agenda

- ☐ Background
- ☐ Aspects in Control Register
- ☐ Aspects in MSR Register
- ☐ Aspects in memory virtualization
- ☐ Aspects in IO
- ☐ Summary

Agenda

☒ Background

☐ Aspects in Control Register

☐ Aspects in MSR Register

☐ Aspects in memory virtualization

☐ Aspects in IO

☐ Summary

Three different context

Virtualization concept 101



Dong, Eddie

To: OTC PRC AutoHV

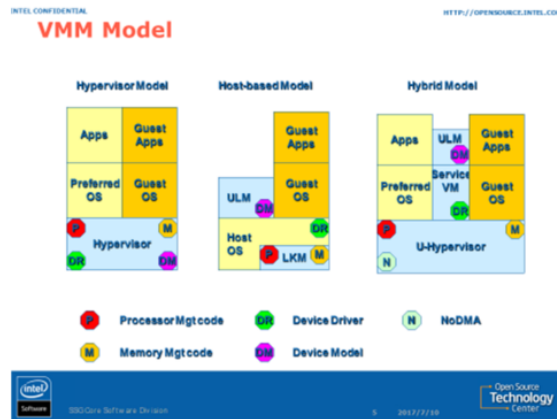
Cc: OTC IAH architects; OTC Cicada Dev; OTC PRC VGT; Dong, Eddie

You forwarded this message on 7/23/2018 5:59 PM.

Outside Intel, VMware & academia people category VMMs into 2: Type-1 & Type-2. Type-1 means Hypervisor Model & Hybrid Model. Type-2 means Host-based model.

Reply Reply All Forward

Tue 7/11/2017 4:15 A



One of the key concept a VM developer must be aware is the 3 different context: Virtual context & Machine Context or physical context, and shadow context. Take CPU as example here. We usually use vCPU, pCPU, vReg, pReg & shadow_Reg to distinguish them. The vXXX context/state is the context/state the VM sees (when the VM is executing). For example, if a guest executes an instruction, MOV AX, CR0 (AX=CR0). The context of CR0 the guest instruction get is the vContext, i.e., vCR0. While, in the real CPU side, the CPU CR0 is different. In this case, we usually call it the shadow context, i.e. at that time, the CR0 of pCPU is shadow CR0. Some time, we name it as sCR0, or shadow_CR0. Shadow context means the context the physical CPU/machine has when the VM is executing.

Another important situation is that when VM exit happens, and the VMM runs. At that time, the context of CPU/machine is called pContext or host Context – this is usually different with shadow context. This is the state/context the real CPU/machine has when the CPU is executing VMM.

Three different context

- Physical context
- Shadow context
- Virtual context

系统虚拟化

第一类原因是访问了特权资源，对 CR 和 MSR 寄存器的访问都属于这一类。

对于此类 VM-Exit，VMM 通过特权资源的虚拟化来解决。特权资源虚拟化的要点在于解决客户机与 VMM 在特权资源控制权的矛盾。即客户机认为自己完全拥有特权资源，可以

Do Not Copy Or Share

自由读写，而特权资源的实际拥有者是 VMM，不能允许客户机自由读写。VMM 通过引入“虚拟特权资源”和“影子特权资源”来解决这个矛盾。“虚拟特权资源”是客户机所看到的特权资源，VMM 允许客户机自由的读写。“影子特权资源”是客户机运行时特权资源真正的值，通常是 VMM 在“虚拟特权资源”的基础上经过处理得到的，因此称其为“影子”。

图 5-8 以特权寄存器为例，展示了特权寄存器的虚拟化过程。当 VCPU 读特权寄存器时，VMM 将“虚拟寄存器”的值返回。比如对于“MOV EAX, CR0”指令，VMM 将 Virtual CR0 的值赋给 EAX，然后 VM-Entry 返回。当 VCPU 写特权寄存器时，VMM 首先将值写入“虚拟寄存器”，然后根据“虚拟寄存器”的值以及虚拟化策略来更新“影子寄存器”，最后将“影子寄存器”的值应用到 VCPU 上，将值写入 VMCS “客户机状态域”的对应字段并且 VM-Entry 返回。这里的虚拟化策略是因特权虚拟器而异的，比如对于如下指令

Agenda

- ☐ Background
- ☒ Aspects in Control Register
- ☐ Aspects in MSR Register
- ☐ Aspects in memory virtualization
- ☐ Aspects in IO
- ☐ Summary

CR4 register

- **Physical context** (the physical value in root mode)

vmcs VMX_HOST_CR4 field.

- **Shadow context** (the physical value in non-root mode)

For the bit **owned by guest**, the value is derived from guest. For the bit **owned by host**, the value is derived from vmcs VMX_GUEST_CR4 field.

- **Virtual context**

For the bit **owned by guest**, the value is derived from guest. For the bit **owned by host**, the value is derived from vmcs VMX_CR4_READ_SHADOW field

Rethinking

- Why CR4.VMXE is **TRAP_AND_EMULATE_BITS**^[1]

23.7 ENABLING AND ENTERING VMX OPERATION

Before system software can enter VMX operation, it enables VMX by setting CR4.VMXE[bit 13] = 1. VMX operation is then entered by executing the VMXON instruction. VMXON causes an invalid-opcode exception (#UD) if executed with CR4.VMXE = 0. Once in VMX operation, it is not possible to clear CR4.VMXE (see Section 23.8). System software leaves VMX operation by executing the VMXOFF instruction. CR4.VMXE can be cleared outside of VMX operation after executing of VMXOFF.

Need to ensure vmxe bit is set in shadow context

^[1] [ACRN mainline virtual_cr.c](#)

Rethinking

- Why CR0.PG is **TRAP_AND_PASSTHRU_BITS** [1]

```
if ((cr0_changed_bits & CR0_PG) != 0UL) {  
    /* PG bit changes */  
    if ((effective_cr0 & CR0_PG) != 0UL) {  
        /* Enable paging */  
        if ((vcpu_get_efer(vcpu) & MSR_IA32_EFER_LME_BIT) != 0UL) {  
            /* Enable long mode */  
            pr_dbg("VMM: Enable long mode");  
            entry_ctrls = exec_vmread32(VMX_ENTRY_CONTROLS);  
            entry_ctrls |= VMX_ENTRY_CTLS_IA32E_MODE;  
            exec_vmwrite32(VMX_ENTRY_CONTROLS, entry_ctrls);  
  
            vcpu_set_efer(vcpu, vcpu_get_efer(vcpu) | MSR_IA32_EFER_LMA_BIT);  
        }  
    }  
}
```

In non-root mode, hardware couldn't update vmcs
VMX_ENTRY_CONTROLS field.
Update this field in root mode.

[1] [ACRN mainline virtual_cr.c](#)

VMXON instruction

Protected Mode Exceptions

#GP(0)	<p>If executed outside VMX operation with CPL>0 or with invalid CR0 or CR4 fixed bits.</p> <p>If executed in A20M mode.</p> <p>If the memory source operand effective address is outside the CS, DS, ES, FS, or GS segment limit.</p> <p>If the DS, ES, FS, or GS register contains an unusable segment.</p> <p>If the source operand is located in an execute-only code segment.</p> <p>If the value of the IA32_FEATURE_CONTROL MSR does not support entry to VMX operation in the current processor mode.</p>
#PF(fault-code)	<p>If a page fault occurs in accessing the memory source operand.</p>
#SS(0)	<p>If the memory source operand effective address is outside the SS segment limit.</p> <p>If the SS register contains an unusable segment.</p>
#UD	<p>If operand is a register.</p> <p>If executed with CR4.VMXE = 0.</p>

That means the physical value of CR4.VMXE.
It's **shadow context** in non-root mode.

Conclusion

- If one bit has restriction in VMX operation or needs to do some operations in root mode, It's better to trap(owned by host) this bit.

Agenda

- ☐ Background
- ☐ Aspects in Control Register
- ☒ **Aspects in MSR Register**
- ☐ Aspects in memory virtualization
- ☐ Aspects in IO
- ☐ Summary

Hardware-assisted save and restore MSR

- VMCS field MSR
- MSR area
 - A VMM may specify lists of MSRs to be stored and loaded on VM exits
 - A VMM may specify a list of MSRs to be loaded on VM entries

We can get physical context and shadow context from above information block.

virtual context may diff with **shadow context**,
it's MSR specific, **it depends on HV.**

VMCS field MSR(MSR_IA32_EFER)

- VMX_HOST_IA32_EFER_FULL
- VMX_GUEST_IA32_EFER_FULL
- HV has ensured shadow context equals to virtual context, so don't need to intercept RDMSR

MSR area(MSR_IA32_TSC_AUX)

- In ACRN, virtual context equals to shadow context

Non hardware-assisted save and restore MSR

- shadow context = physical context
- MSR_IA32_EXT_XAPICID
- MSR_IA32_TIME_STAMP_COUNTER
- MSR_IA32_EXT_APIC_VERSION

virtual context may diff with **shadow context**,
it's MSR specific.

MSR_IA32_EXT_XAPICID

- For isolation, shadow context diff with virtual context
- Need to intercept RDMSR
- Virtual context value is set by HV

MSR_IA32_TIME_STAMP_COUNTER

- For isolation, shadow context diff with virtual context
- Don't need to intercept RDMSR

25.3 CHANGES TO INSTRUCTION BEHAVIOR IN VMX NON-ROOT OPERATION

RDMSR. Section 25.1.3 identifies when executions of the RDMSR instruction cause VM exits. If such an execution causes neither a fault due to CPL > 0 nor a VM exit, the instruction's behavior may be modified for certain values of ECX:

- If ECX contains 10H (indicating the IA32_TIME_STAMP_COUNTER MSR), the value returned by the instruction is determined by the setting of the "use TSC offsetting" VM-execution control:
 - If the control is 0, RDMSR operates normally, loading EAX:EDX with the value of the IA32_TIME_STAMP_COUNTER MSR.
 - If the control is 1, the value returned is determined by the setting of the "use TSC scaling" VM-execution control:

A logical processor uses PAE paging if CR0.PG = 1, CR4.PAE = 1 and IA32_EFER.LMA = 0. See Section 4.4 in the *Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3A*.

i-8 Vol. 3C

VMX NON-ROOT OPERATION

- If the control is 0, RDMSR loads EAX:EDX with the sum of the value of the IA32_TIME_STAMP_COUNTER MSR and the value of the TSC offset.
- If the control is 1, RDMSR first computes the product of the value of the IA32_TIME_STAMP_COUNTER MSR and the value of the TSC multiplier. It then shifts the value of the product right 48 bits and loads EAX:EDX with the sum of that shifted value and the value of the TSC offset.

MSR_IA32_EXT_APIC_VERSION

- shadow context = virtual context

Agenda

- ☐ Background
- ☐ Aspects in Control Register
- ☐ Aspects in MSR Register
- ☒ **Aspects in memory virtualization**
- ☐ Aspects in IO
- ☐ Summary

EPT

- **Shadow context**

HPA

- **Virtual context**

GPA

- **Physical context**

HPA, but may need the relationship between GPA, HVA and HPA to build EPT table

Agenda

- ☐ Background
- ☐ Aspects in Control Register
- ☐ Aspects in MSR Register
- ☐ Aspects in memory virtualization
- ☒ Aspects in IO
- ☐ Summary

MSI interrupt W/O interrupt remapping

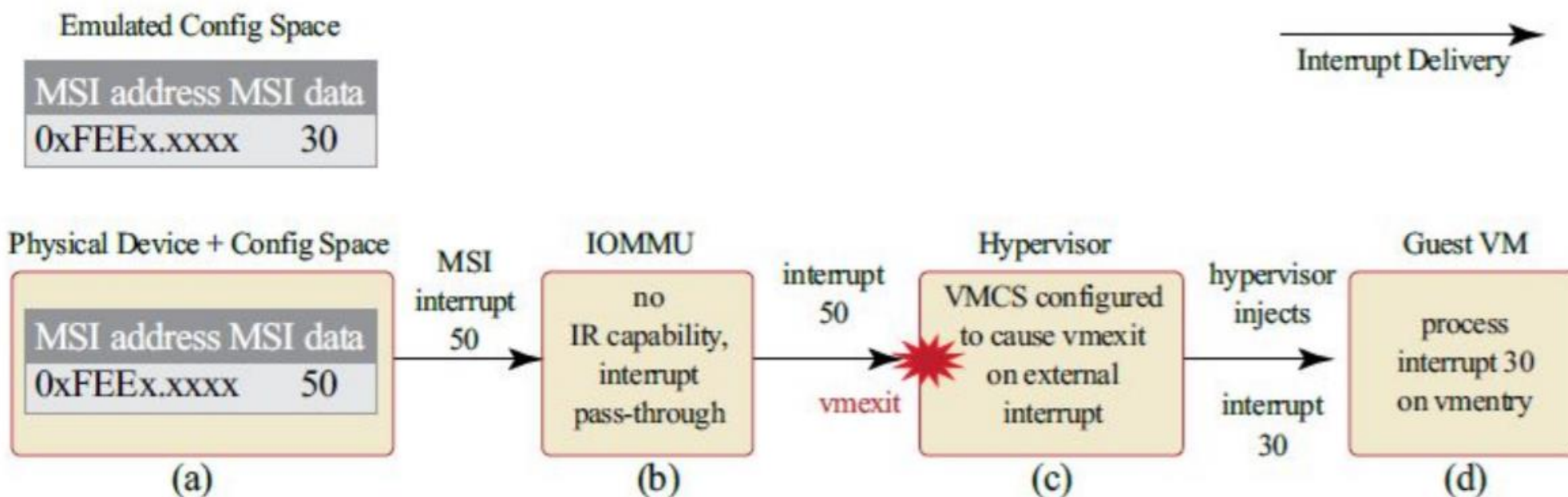


Figure 6.16: MSI interrupt delivery without interrupt remapping support.

MSI interrupt With interrupt remapping

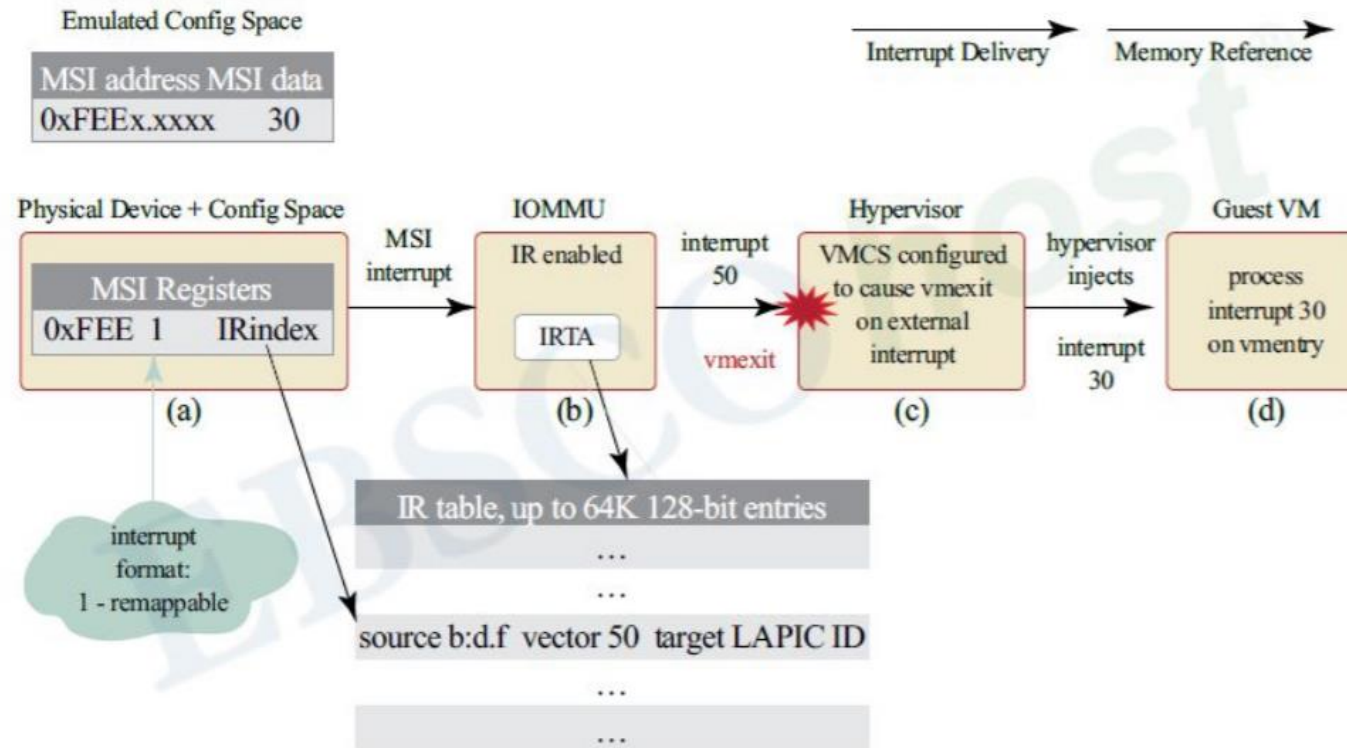


Figure 6.17: MSI interrupt delivery with interrupt remapping support. (IRindex is denoted “interrupt_index” in the VT-d specification.)

MSI interrupt

- In the previous two cases, shadow context equals to physical context
- virtual context differs with shadow context

Agenda

- ☐ Background
- ☐ Aspects in Control Register
- ☐ Aspects in MSR Register
- ☐ Aspects in memory virtualization
- ☐ Aspects in IO
- ☒ Summary

Summary

- Must be aware of the 3 different context
- Know what is done by hardware, what is done by software
- Know what is done in root mode, what is done in non-root mode