

# Cloud Operating Systems

PIC and PIT

**Andreas Kogler, Fabian Rauscher, Daniel Gruss**

2023-03-27









- Classic hardware interactions



- Classic hardware interactions
- What is the PIC?



- Classic hardware interactions
- What is the PIC?
- What is the PIT?



- Classic hardware interactions
- What is the PIC?
- What is the PIT?
- Why should *your* **HV** care?





- Classic hardware interactions
- What is the PIC?
- What is the PIT?
- Why should *your* **HV** care?
- How to **virtualize** the PIC and PIT

**Hardware Interactions?**









- draw a **character** on the console?



- draw a **character** on the console?
- receive a **keyboard** press?



- draw a **character** on the console?
- receive a **keyboard** press?
- receive a regular *heartbeat*?











- **frame buffer** at physical address `0xB8000`



- **frame buffer** at physical address `0xB8000`
- how to pass the buffer to the guest?



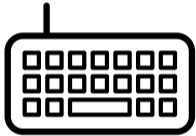
- **frame buffer** at physical address `0xB8000`
- how to pass the buffer to the guest?
  - share it via the EPT?

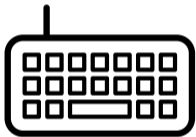


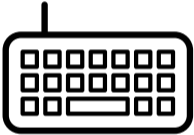
- **frame buffer** at physical address `0xB8000`
- how to pass the buffer to the guest?
  - share it via the EPT?
  - copy it, but **when**?



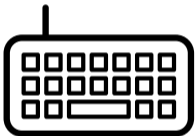




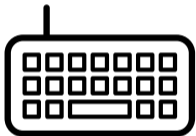




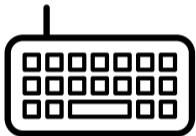
- implemented via **interrupts** and **ports**



- implemented via **interrupts** and **ports**
- check if data is available on port `0x64`



- implemented via **interrupts** and **ports**
- check if data is available on port `0x64`
- get the current *scancode* key value from port `0x60`



- implemented via **interrupts** and **ports**
- check if data is available on port `0x64`
- get the current *scancode* key value from port `0x60`
- but how to share with a guest?











- How and when to **preempt** threads?



- How and when to **preempt** threads?
- Timer waking up the scheduler



- How and when to **preempt** threads?
- Timer waking up the scheduler
- The scheduler decides what to do



- How and when to **preempt** threads?
- Timer waking up the scheduler
- The scheduler decides what to do
- But how to pass this heartbeat to the guest?

# Programmable Interrupt Controller











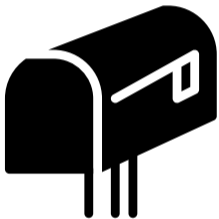
- The PIC makes x86 **interrupt driven**



- The PIC makes x86 **interrupt driven**
- Manages *hardware* interrupts



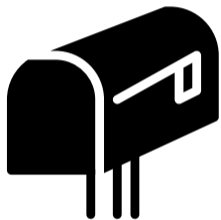
- The PIC makes x86 **interrupt driven**
- Manages *hardware* interrupts
- Wires *hardware* interrupts to *system* interrupts



- The PIC makes x86 **interrupt driven**
- Manages *hardware* interrupts
- Wires *hardware* interrupts to *system* interrupts
- More modern systems use the successor: **APIC**



- The PIC makes x86 **interrupt driven**
- Manages *hardware* interrupts
- Wires *hardware* interrupts to *system* interrupts
- More modern systems use the successor: **APIC**
- SWEB uses the PIC



- The PIC makes x86 **interrupt driven**
- Manages *hardware* interrupts
- Wires *hardware* interrupts to *system* interrupts
- More modern systems use the successor: **APIC**
- SWEB uses the PIC
- Checkout: [https://wiki.osdev.org/8259\\_PIC](https://wiki.osdev.org/8259_PIC)











- Features 16 interrupts



- Features 16 interrupts
- Notation: **Port[P] = D**



- Features 16 interrupts
- Notation: **Port[P] = D**
- distributed to *Parent* (**P=0x20**) and *Child* (**P=0xA0**) PIC



- Features 16 interrupts
- Notation: **Port[P] = D**
- distributed to *Parent* (**P=0x20**) and *Child* (**P=0xA0**) PIC
- Data (**P+1**) and command (**P+0**) **ports**



- Features 16 interrupts
- Notation: **Port[P] = D**
- distributed to *Parent* (**P=0x20**) and *Child* (**P=0xA0**) PIC
- Data (**P+1**) and command (**P+0**) **ports**
- Most known command: End-Of-Interrupt (EOI) **D=0x20**











- SWEB only uses the PIC in *one* configuration



- SWEB only uses the PIC in *one* configuration
- Initial write of the config ports with **D=0x11**



- SWEB only uses the PIC in *one* configuration
- Initial write of the config ports with **D=0x11**
- Followed by **three** writes to the data *port*



- SWEB only uses the PIC in *one* configuration
- Initial write of the config ports with **D=0x11**
- Followed by **three** writes to the data *port*
  - First, the **offset** to translate hw IRQ to vector numbers



- SWEB only uses the PIC in *one* configuration
- Initial write of the config ports with **D=0x11**
- Followed by **three** writes to the data *port*
  - First, the **offset** to translate hw IRQ to vector numbers
- Further writes only change the interrupt **mask**



- SWEB only uses the PIC in *one* configuration
- Initial write of the config ports with **D=0x11**
- Followed by **three** writes to the data *port*
  - First, the **offset** to translate hw IRQ to vector numbers
- Further writes only change the interrupt **mask**
- Important hardware IRQs:





- SWEB only uses the PIC in *one* configuration
- Initial write of the config ports with **D=0x11**
- Followed by **three** writes to the data *port*
  - First, the **offset** to translate hw IRQ to vector numbers
- Further writes only change the interrupt **mask**
- Important hardware IRQs:
  - **0**: Timer interrupt



- SWEB only uses the PIC in *one* configuration
- Initial write of the config ports with **D=0x11**
- Followed by **three** writes to the data *port*
  - First, the **offset** to translate hw IRQ to vector numbers
- Further writes only change the interrupt **mask**
- Important hardware IRQs:
  - **0**: Timer interrupt
  - **1**: Keyboard interrupt

# Programmable Interval Timer









- Basically a programmable oscillator, with selectable:



- Basically a programmable oscillator, with selectable:
  - frequency





- Basically a programmable oscillator, with selectable:
  - frequency
  - divisor



- Basically a programmable oscillator, with selectable:
  - frequency
  - divisor
  - modes



- Basically a programmable oscillator, with selectable:
  - frequency
  - divisor
  - modes
- SWEB uses channel 0 (**P=0x40**)

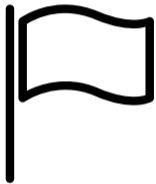


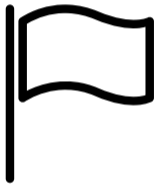
- Basically a programmable oscillator, with selectable:
  - frequency
  - divisor
  - modes
- SWEB uses channel 0 (**P=0x40**)
- Including command port (**P=0x43**) with (**D=0x36**)



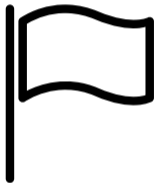
- Basically a programmable oscillator, with selectable:
  - frequency
  - divisor
  - modes
- SWEB uses channel 0 (**P=0x40**)
- Including command port (**P=0x43**) with (**D=0x36**)
- Checkout:  
[https://wiki.osdev.org/Programmable\\_Interval\\_Timer](https://wiki.osdev.org/Programmable_Interval_Timer)







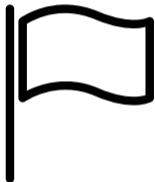




- Needed by the **scheduler** to change threads



- Needed by the **scheduler** to change threads
- Channel 0 counts with 18.2065 Hz



- Needed by the **scheduler** to change threads
- Channel 0 counts with 18.2065 Hz
- Giving an **interrupt** every  $\approx 54$  ms



- Needed by the **scheduler** to change threads
- Channel 0 counts with 18.2065 Hz
- Giving an **interrupt** every  $\approx 54$  ms
- If enabled in the **PIC**

**How to virtualize PIC and PIT?**

1. Check pending interrupts  
→ Keyboard, Keyboard

Int. Window Exiting	0
Int. Info.	0
rflags.IF	0
pending IRQs	Keyboard

1. Check pending interrupts  
→ Keyboard, Keyboard
2. Check guest IF flag → 0

Int. Window Exiting	0
Int. Info.	0
rflags.IF	0
pending IRQs	Keyboard

1. Check pending interrupts  
→ Keyboard, Keyboard
2. Check guest IF flag → 0
3. Set interrupt-window exiting

Int. Window Exiting	1
Int. Info.	0
rflags.IF	0
pending IRQs	Keyboard



1. Check pending interrupts  
→ Keyboard, Keyboard
2. Check guest IF flag → 0
3. Set interrupt-window exiting
4. Continue guest

Int. Window Exiting	1
Int. Info.	0
rflags.IF	0
pending IRQs	Keyboard

```
push  rbp
mov   rbp, rsp
sti
nop
pop   rbp
ret
```

Int. Window Exiting	1
Int. Info.	0
rflags.IF	0
pending IRQs	Keyboard

```
push    rbp
mov     rbp, rsp
sti
nop
pop     rbp
ret
```

Int. Window Exiting	1
Int. Info.	0
rflags.IF	0
pending IRQs	Keyboard

```
push    rbp
mov     rbp, rsp
sti
nop
pop     rbp
ret
```

Int. Window Exiting	1
Int. Info.	0
rflags.IF	0
pending IRQs	Keyboard

```
push    rbp
mov     rbp, rsp
sti
nop
pop     rbp
ret
```

Int. Window Exiting	1
Int. Info.	0
rflags.IF	1
pending IRQs	Keyboard

```
push    rbp
mov     rbp, rsp
sti
nop
pop     rbp
ret
```



Int. Window Exiting	1
Int. Info.	0
rflags.IF	1
pending IRQs	Keyboard

1. VM Exit → VMM takes control

Int. Window Exiting	1
Int. Info.	0
rflags.IF	1
pending IRQs	Keyboard

1. VM Exit → VMM takes control
2. Check VM Exit reason  
→ interrupt window

Int. Window Exiting	1
Int. Info.	0
rflags.IF	1
pending IRQs	Keyboard



1. VM Exit → VMM takes control
2. Check VM Exit reason
  - interrupt window
  - 2.1 Check pending interrupts
    - Keyboard, Keyboard

Int. Window Exiting	1
Int. Info.	0
rflags.IF	1
pending IRQs	Keyboard

1. VM Exit → VMM takes control
2. Check VM Exit reason
  - interrupt window
  - 2.1 Check pending interrupts
    - Keyboard, Keyboard
  - 2.2 Set interruption information

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	1
pending IRQs	Keyboard

1. VM Exit → VMM takes control
2. Check VM Exit reason
  - interrupt window
    - 2.1 Check pending interrupts
      - Keyboard, Keyboard
    - 2.2 Set interruption information
    - 2.3 Continue guest

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	1
pending IRQs	

```
push    rbp
mov     rbp, rsp
sti
nop
pop     rbp
ret
```

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	1
pending IRQs	

```
push    rbp
mov     rbp, rsp
sti
nop
pop     rbp
ret
```



Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	1
pending IRQs	

```
<irq_handler>
```

```
in EAX, 0x64
```

```
iret
```

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	0
pending IRQs	

<irq handler>

in **EAX**, 0x64

iret



Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	0
pending IRQs	

1. IN and OUT always use EAX, AL, AX

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	0
pending IRQs	



1. IN and OUT always use EAX, AL, AX
2. Hypervisor injects value
3. Hypervisor advances RIP

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	0
pending IRQs	

1. IN and OUT always use EAX, AL, AX
2. Hypervisor injects value
3. Hypervisor advances RIP
4. Continue Guest

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	0
pending IRQs	

1. IN and OUT always use EAX, AL, AX
2. Hypervisor injects value
3. Hypervisor advances RIP
4. Continue Guest

Int. Window Exiting	0
Int. Info.	0x80000021
rflags.IF	0
pending IRQs	









- Examine VMExits based on IO port operations



- Examine VMExits based on IO port operations
- Emulate PIC and PIT for their ports





- Examine VMExits based on IO port operations
- Emulate PIC and PIT for their ports
- You don't need to implement all features



- Examine VMExits based on IO port operations
- Emulate PIC and PIT for their ports
- You don't need to implement all features
  - only what SWEB needs



- Examine VMExits based on IO port operations
- Emulate PIC and PIT for their ports
- You don't need to implement all features
  - only what SWEB needs
  - check SWEB's code if unsure









- Emulate the ports for the configuration:



- Emulate the ports for the configuration:
  - the **offset** mapping HW IRQ to interrupt vectors





- Emulate the ports for the configuration:
  - the **offset** mapping HW IRQ to interrupt vectors
  - the interrupt mask enabling and disabling interrupts



- Emulate the ports for the configuration:
  - the **offset** mapping HW IRQ to interrupt vectors
  - the interrupt mask enabling and disabling interrupts
  - ...



- Emulate the ports for the configuration:
  - the **offset** mapping HW IRQ to interrupt vectors
  - the interrupt mask enabling and disabling interrupts
  - ...
- Forward only enabled interrupts to the guest



- Emulate the ports for the configuration:
  - the **offset** mapping HW IRQ to interrupt vectors
  - the interrupt mask enabling and disabling interrupts
  - ...
- Forward only enabled interrupts to the guest
- Manage pending interrupts in the HV



- Emulate the ports for the configuration:
  - the **offset** mapping HW IRQ to interrupt vectors
  - the interrupt mask enabling and disabling interrupts
  - ...
- Forward only enabled interrupts to the guest
- Manage pending interrupts in the HV
- Think about interrupts during interrupts (hint: EOI)



- Emulate the ports for the configuration:
  - the **offset** mapping HW IRQ to interrupt vectors
  - the interrupt mask enabling and disabling interrupts
  - ...
- Forward only enabled interrupts to the guest
- Manage pending interrupts in the HV
- Think about interrupts during interrupts (hint: EOI)
- Goal: Be able to forward keyboard scancodes to the guest











- Emulate the ports similar to the PIC



- Emulate the ports similar to the PIC
- Generate timer interrupts if enabled in the PIC



- Emulate the ports similar to the PIC
- Generate timer interrupts if enabled in the PIC
- Goal: Generate timer interrupts for the guest's scheduler









- Register: CPU\_BASED\_VM\_EXEC\_CONTROL





- Register: CPU\_BASED\_VM\_EXEC\_CONTROL
  - Bit: CPU\_BASED\_VIRTUAL\_INTR\_PENDING



- Register: CPU\_BASED\_VM\_EXEC\_CONTROL
  - Bit: CPU\_BASED\_VIRTUAL\_INTR\_PENDING
- Register: VM\_ENTRY\_INTR\_INFO\_FIELD



- Register: CPU\_BASED\_VM\_EXEC\_CONTROL
  - Bit: CPU\_BASED\_VIRTUAL\_INTR\_PENDING
- Register: VM\_ENTRY\_INTR\_INFO\_FIELD
- Register: VMX\_PREEMPTION\_TIMER\_VALUE



- Register: CPU\_BASED\_VM\_EXEC\_CONTROL
  - Bit: CPU\_BASED\_VIRTUAL\_INTR\_PENDING
- Register: VM\_ENTRY\_INTR\_INFO\_FIELD
- Register: VMX\_PREEMPTION\_TIMER\_VALUE
- ExitReasons: INT\_WINDOW, IO\_INST and NMI

**Questions?**