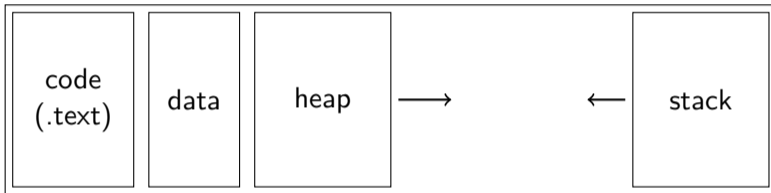


Virtualisation

Johan Montelius

KTH

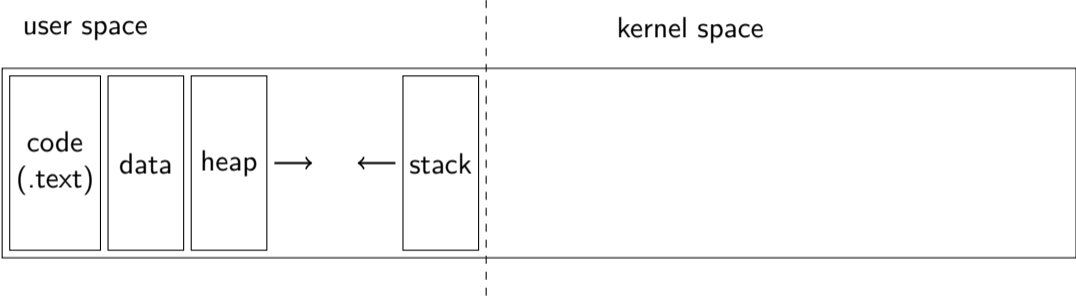
2021



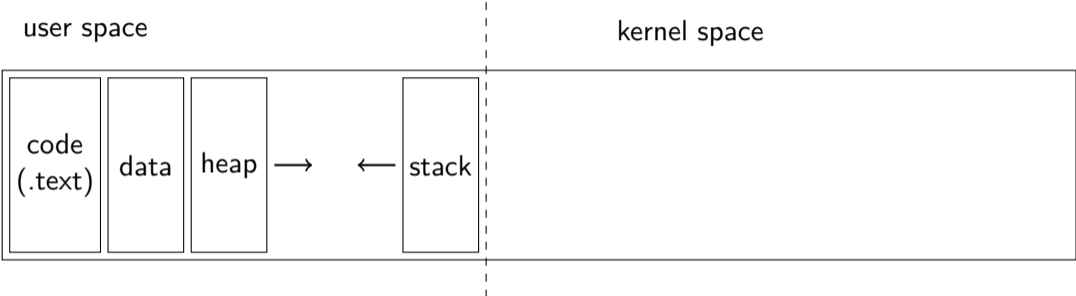
The role of the operating system - provide a virtual environment for a process.



the kernel



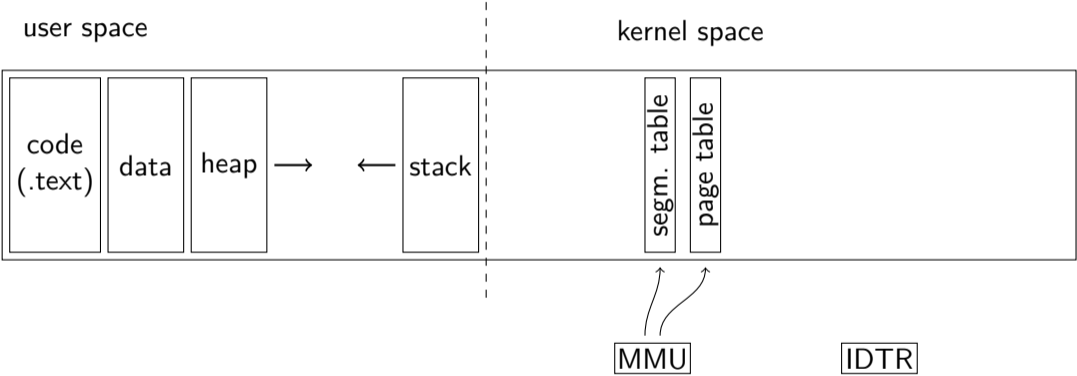
the kernel



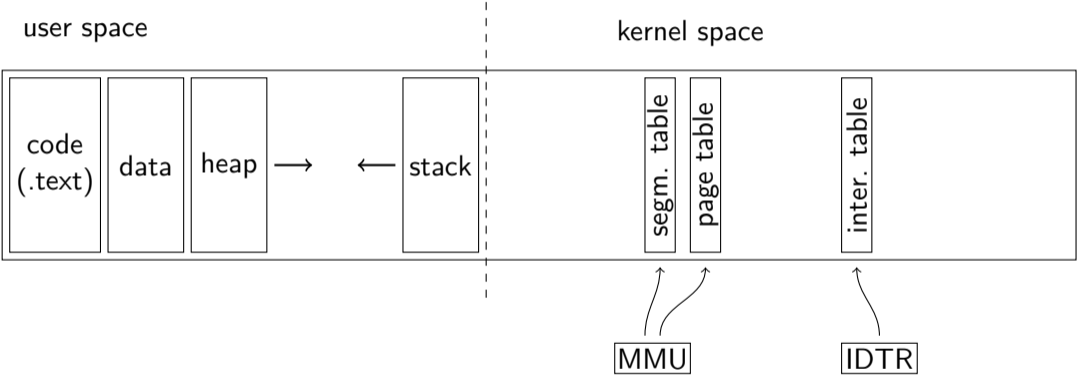
MMU

IDTR

the kernel



the kernel



Who is in control?

Who is in control?

- control the registers of the MMU and you control the virtual address space

Who is in control?

- control the registers of the MMU and you control the virtual address space
- control the IDTR and you control what will happen when we have an interrupt

Who is in control?

- control the registers of the MMU and you control the virtual address space
- control the IDTR and you control what will happen when we have an interrupt
- instructions to set MMU or IDT registers are privileged instructions

Who is in control?

- control the registers of the MMU and you control the virtual address space
- control the IDTR and you control what will happen when we have an interrupt
- instructions to set MMU or IDT registers are privileged instructions

Limited direct execution:

Who is in control?

- control the registers of the MMU and you control the virtual address space
- control the IDTR and you control what will happen when we have an interrupt
- instructions to set MMU or IDT registers are privileged instructions

Limited direct execution:

- only work with mapped memory in user space,

Who is in control?

- control the registers of the MMU and you control the virtual address space
- control the IDTR and you control what will happen when we have an interrupt
- instructions to set MMU or IDT registers are privileged instructions

Limited direct execution:

- only work with mapped memory in user space,
- only execute non-privileged instructions,

Who is in control?

- control the registers of the MMU and you control the virtual address space
- control the IDTR and you control what will happen when we have an interrupt
- instructions to set MMU or IDT registers are privileged instructions

Limited direct execution:

- only work with mapped memory in user space,
- only execute non-privileged instructions,
- for a limited amount of time.

Synchronous interrupts - exceptions:

- faults:

Synchronous interrupts - exceptions:

- faults:
 - page fault
 - privilege violation
 - divide by zero, ...

Synchronous interrupts - exceptions:

- faults:
 - page fault
 - privilege violation
 - divide by zero, ...
- programmed exceptions:

Synchronous interrupts - exceptions:

- faults:
 - page fault
 - privilege violation
 - divide by zero, ...
- programmed exceptions:
 - system call (INT 0x80)
 - debug instructions

Synchronous interrupts - exceptions:

- faults:
 - page fault
 - privilege violation
 - divide by zero, ...
- programmed exceptions:
 - system call (INT 0x80)
 - debug instructions

Asynchronous interrupts:

Synchronous interrupts - exceptions:

- faults:
 - page fault
 - privilege violation
 - divide by zero, ...
- programmed exceptions:
 - system call (INT 0x80)
 - debug instructions

Asynchronous interrupts:

- timer interrupt

Synchronous interrupts - exceptions:

- faults:
 - page fault
 - privilege violation
 - divide by zero, ...
- programmed exceptions:
 - system call (INT 0x80)
 - debug instructions

Asynchronous interrupts:

- timer interrupt
- hardware interrupt: I/O complete, ...

Synchronous interrupts - exceptions:

- faults:
 - page fault
 - privilege violation
 - divide by zero, ...
- programmed exceptions:
 - system call (INT 0x80)
 - debug instructions

Asynchronous interrupts:

- timer interrupt
- hardware interrupt: I/O complete, ...

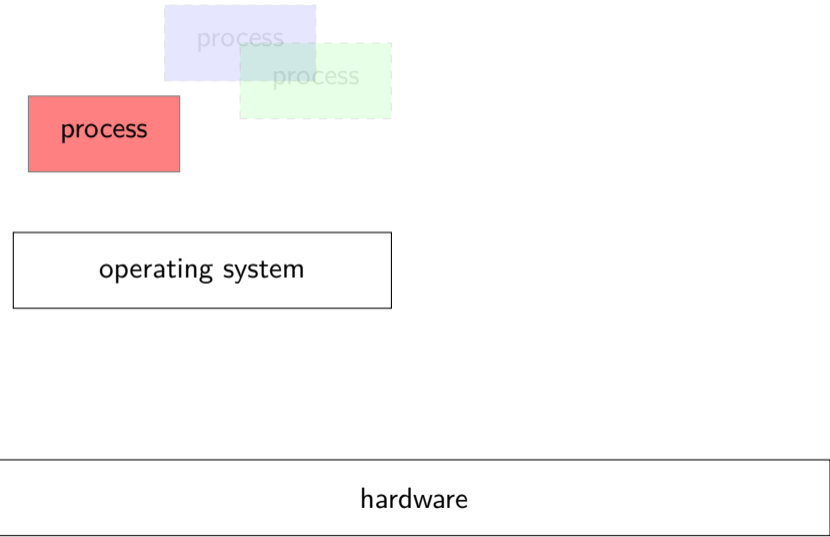
The kernel is interrupt driven.

hardware

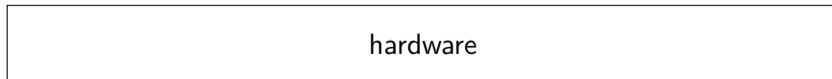
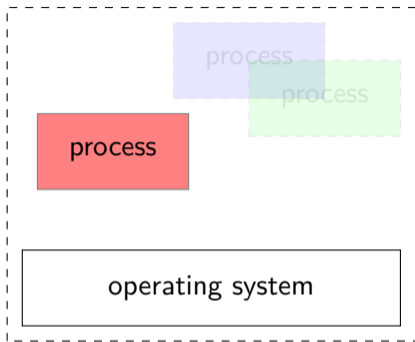
operating system

hardware

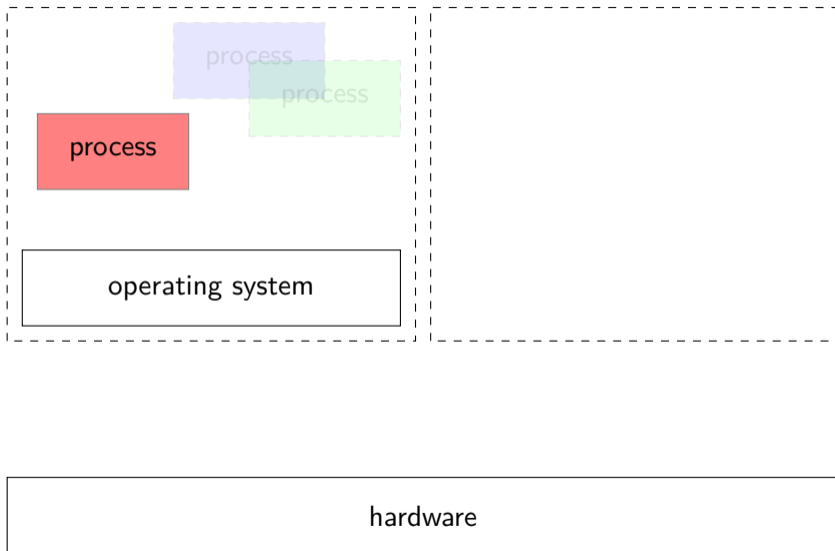
Virtualisation



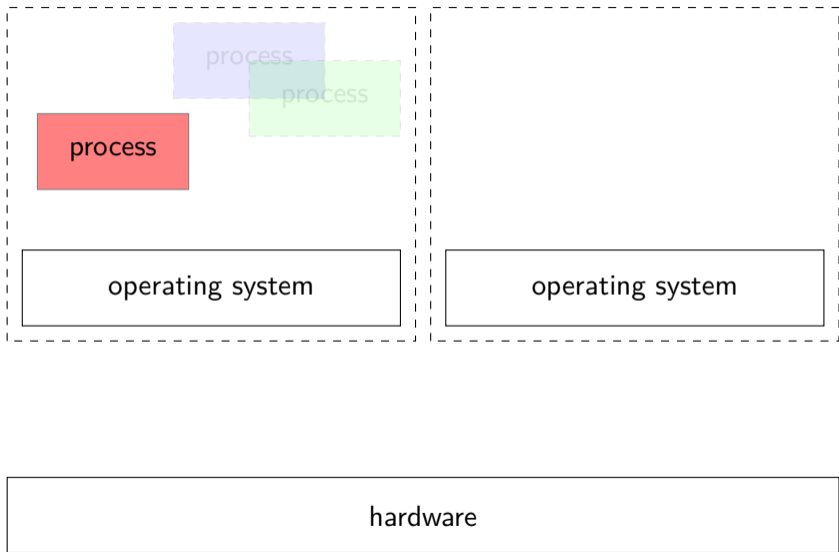
Virtualisation



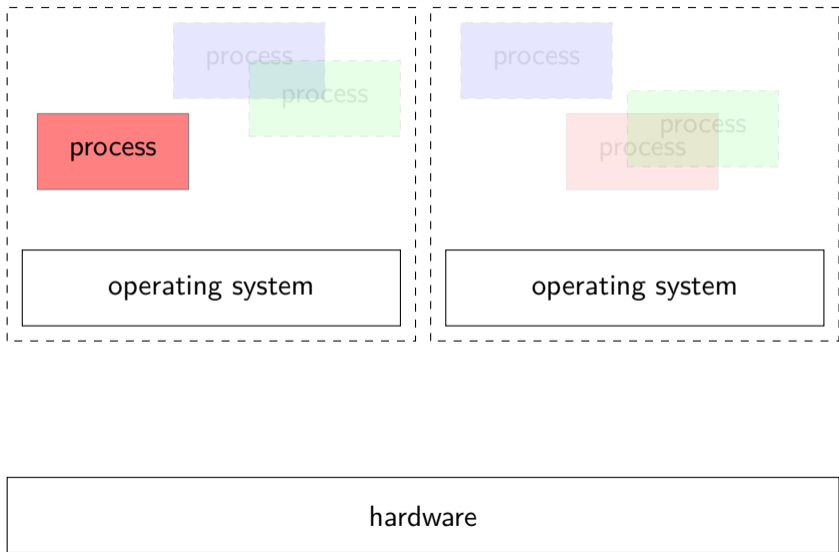
Virtualisation



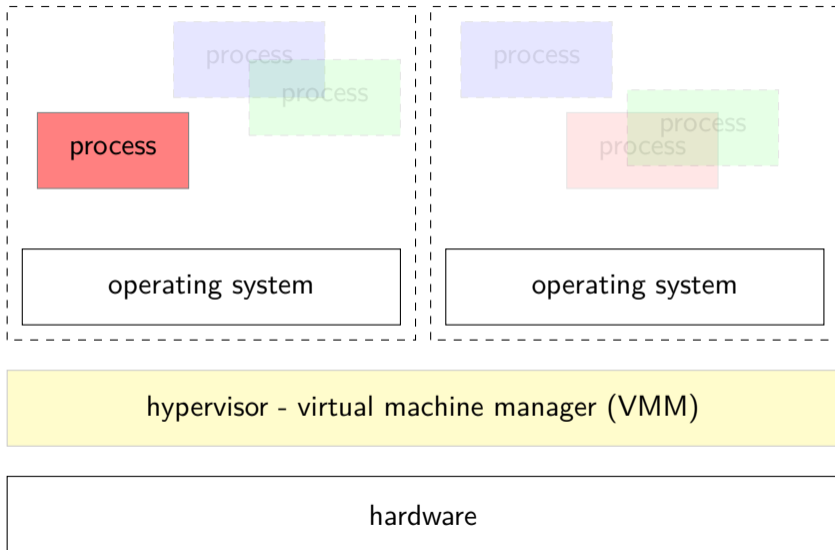
Virtualisation



Virtualisation



Virtualisation



Why?

Why?

Utilisation of hardware.

Utilisation of hardware.

Also provided by a multi-task operating system, what is new?

Utilisation of hardware.

Also provided by a multi-task operating system, what is new?

Applications are completely separated from each other.

Utilisation of hardware.

Also provided by a multi-task operating system, what is new?

Applications are completely separated from each other.

What do two processes in an operating system share?

Utilisation of hardware.

Also provided by a multi-task operating system, what is new?

Applications are completely separated from each other.

What do two processes in an operating system share?

Applications can use different operating systems.

Why?

Utilisation of hardware.

Also provided by a multi-task operating system, what is new?

Applications are completely separated from each other.

What do two processes in an operating system share?

Applications can use different operating systems.

Is this important?

Provide virtualisation of the hardware:

- a virtual cpu, part of the processing power

Provide virtualisation of the hardware:

- a virtual cpu, part of the processing power
- a virtual memory, the illusion of physical memory

Provide virtualisation of the hardware:

- a virtual cpu, part of the processing power
- a virtual memory, the illusion of physical memory

I think we have seen this before.

Provide *limited direct execution* i.e. allow each guest operating system to execute in *user space* and only perform non-privileged operations.

Provide *limited direct execution* i.e. allow each guest operating system to execute in *user space* and only perform non-privileged operations.

What is the first thing an operating system wants to do?

Hypervisor

Guest Operating system

set up IDT

the virtual IDT

Hypervisor

Guest Operating system

set up IDT
pass control to OS

the virtual IDT

Hypervisor

Guest Operating system



the virtual IDT

Hypervisor

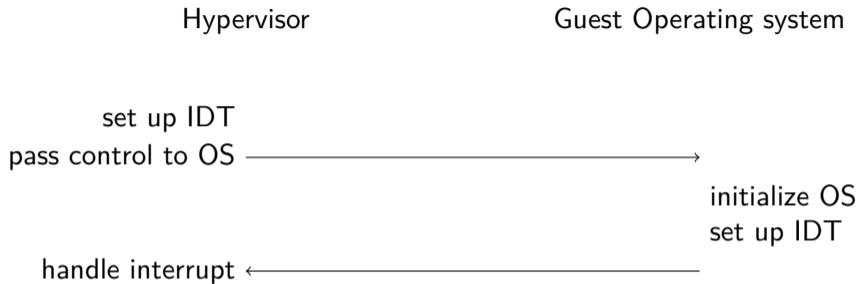
Guest Operating system

set up IDT
pass control to OS



initialize OS
set up IDT

the virtual IDT



the virtual IDT

Hypervisor

Guest Operating system

set up IDT
pass control to OS



initialize OS
set up IDT

handle interrupt
save ref to IDT of OS



the virtual IDT

Hypervisor

Guest Operating system

set up IDT
pass control to OS



initialize OS
set up IDT

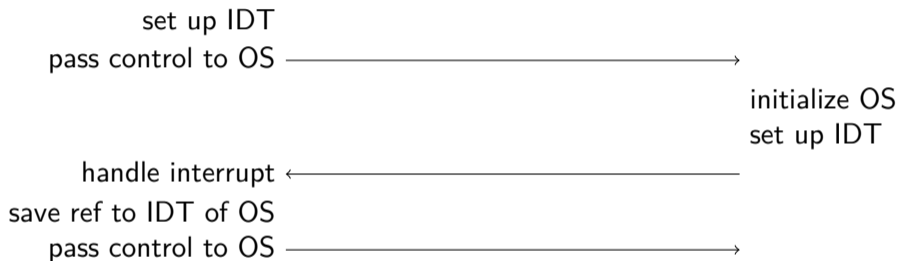
handle interrupt
save ref to IDT of OS
pass control to OS



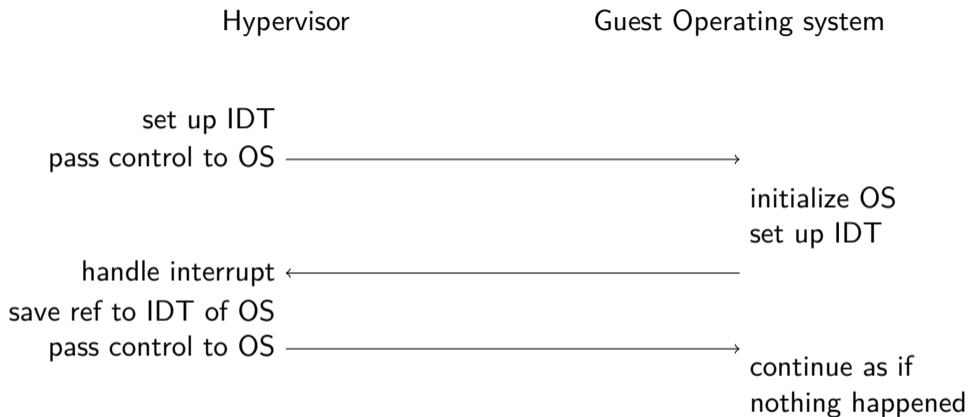
the virtual IDT

Hypervisor

Guest Operating system



the virtual IDT



The operating system is running in non-privileged mode.

a system call

Hypervisor

Guest operating system

Application

running

a system call

Hypervisor

Guest operating system

Application

running
system call

a system call

Hypervisor

Guest operating system

Application

running
system call
INT 0x80

a system call

Hypervisor

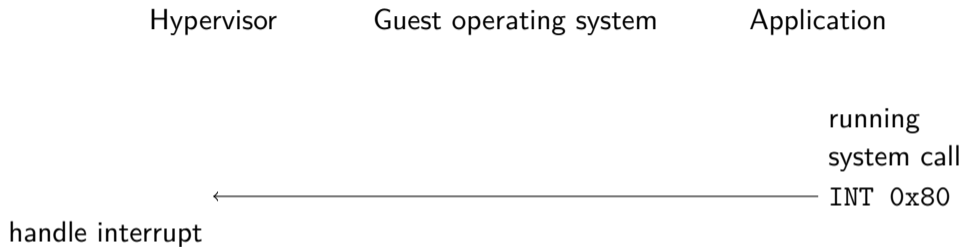
Guest operating system

Application

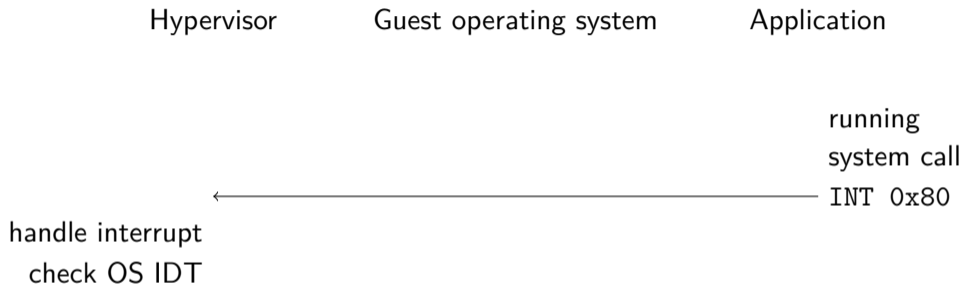
running
system call
INT 0x80



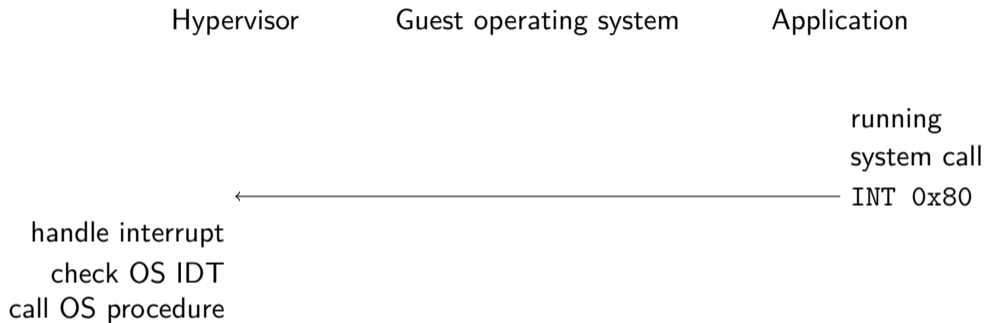
a system call



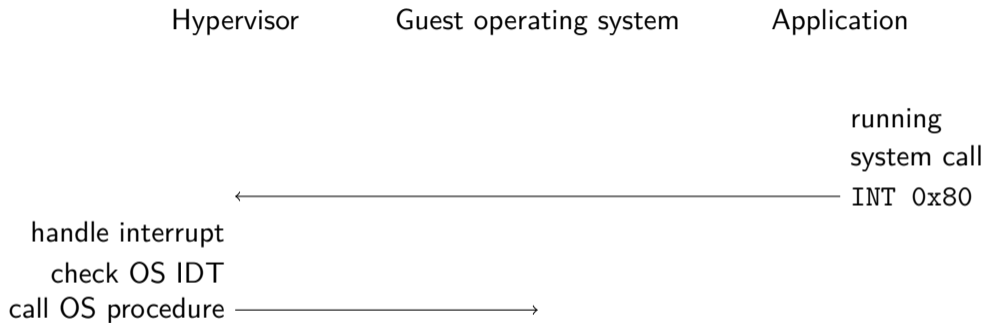
a system call



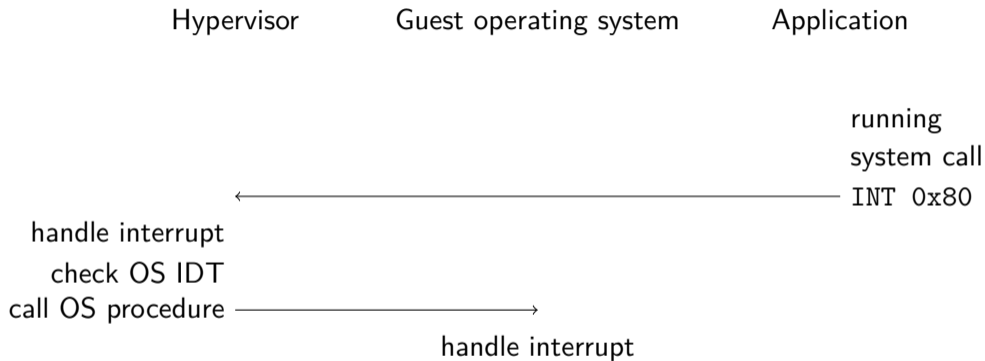
a system call



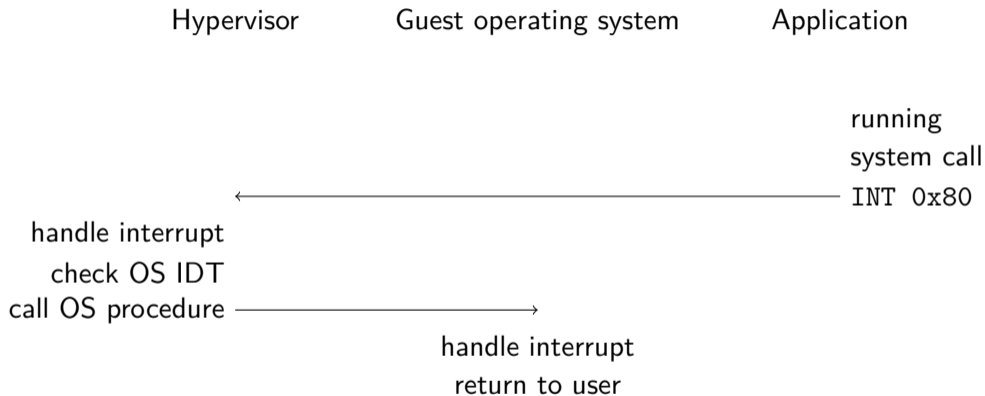
a system call



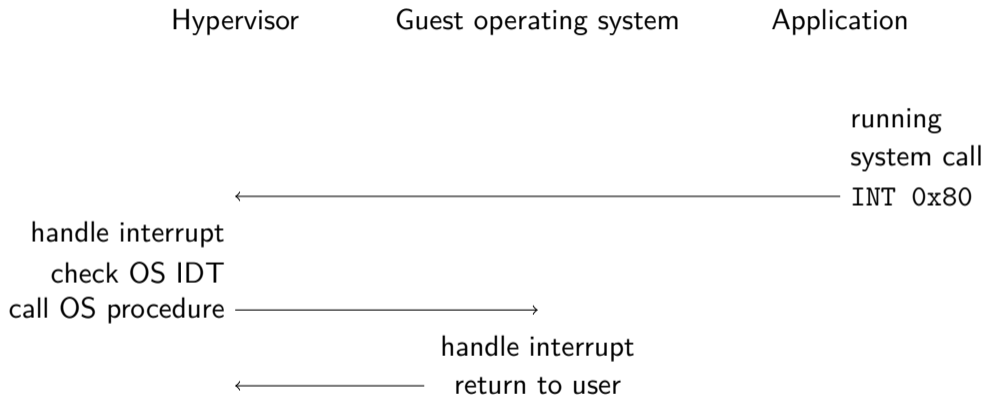
a system call



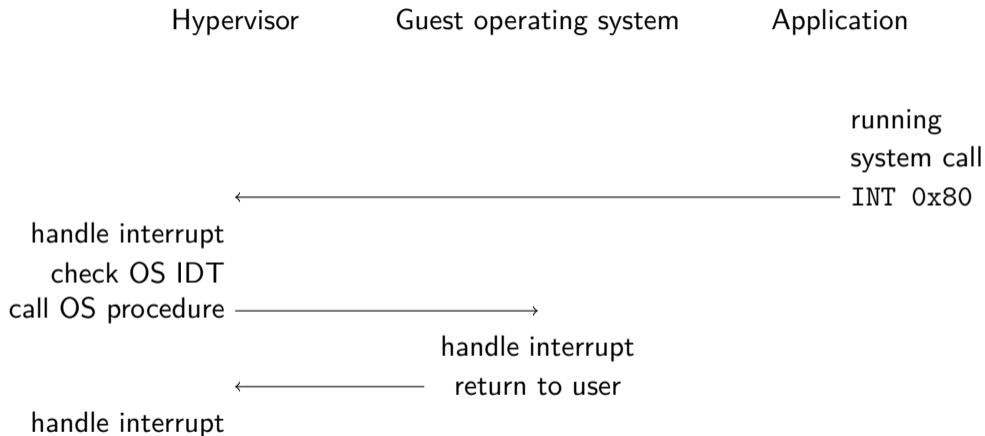
a system call



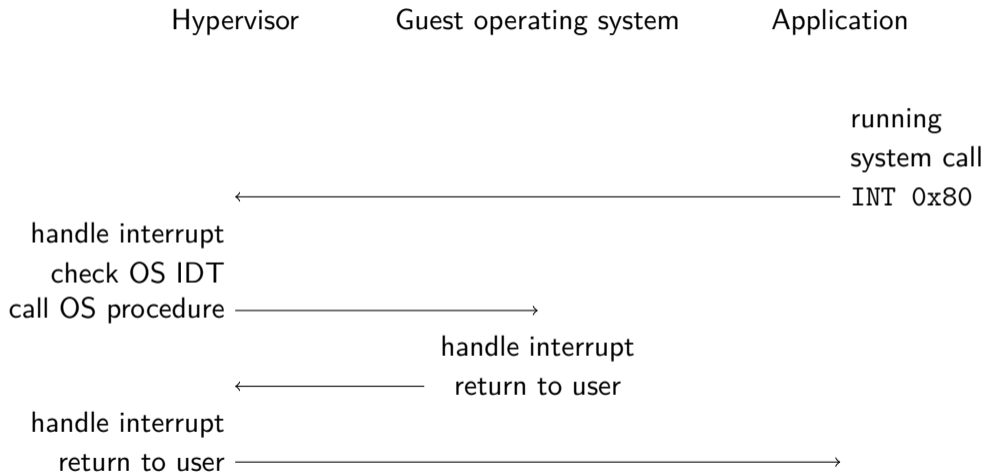
a system call



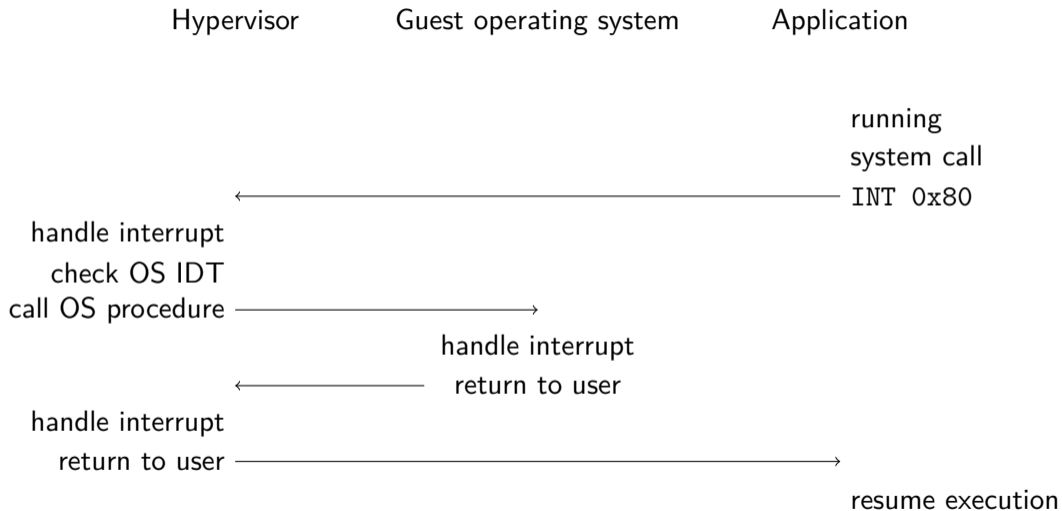
a system call



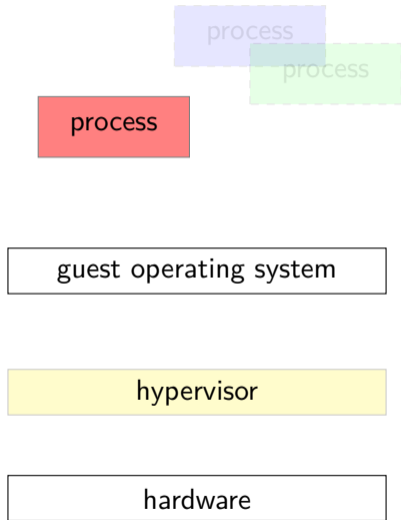
a system call



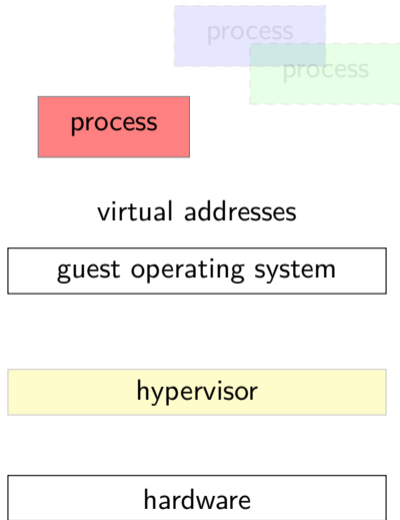
a system call



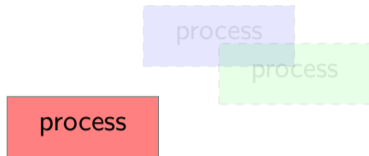
What about virtual memory?



What about virtual memory?



What about virtual memory?



virtual addresses

guest operating system

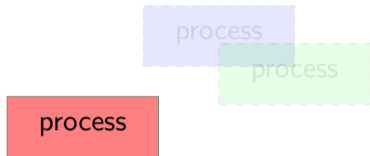
physical addresses

hypervisor

hardware

- regular translation tables

What about virtual memory?



virtual addresses

guest operating system

physical addresses

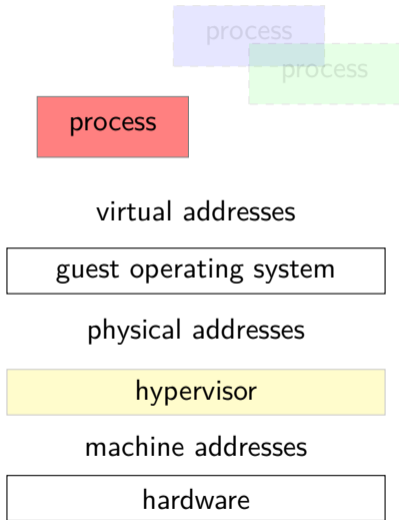
hypervisor

machine addresses

hardware

- regular translation tables
- second level translation

What about virtual memory?



- regular translation tables
- second level translation

This will be expensive!

User process uses virtual addresses that are automatic translated by the hardware (using page table and the MMU) to physical addresses.

User process uses virtual addresses that are automatic translated by the hardware (using page table and the MMU) to physical addresses.

A *page fault* invokes the kernel that, if allowed, maps a missing page and return to the user process.

Hypervisor

Guest operating system

Application

running

Hypervisor

Guest operating system

Application

running
page fault

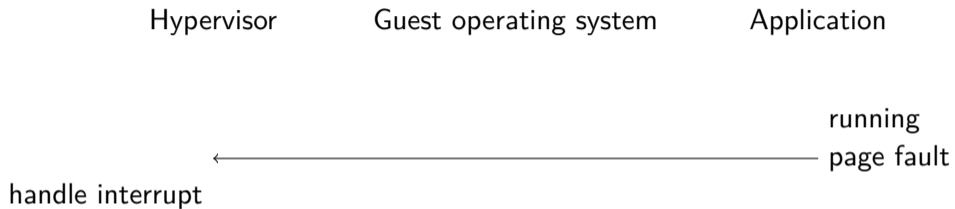
Hypervisor

Guest operating system

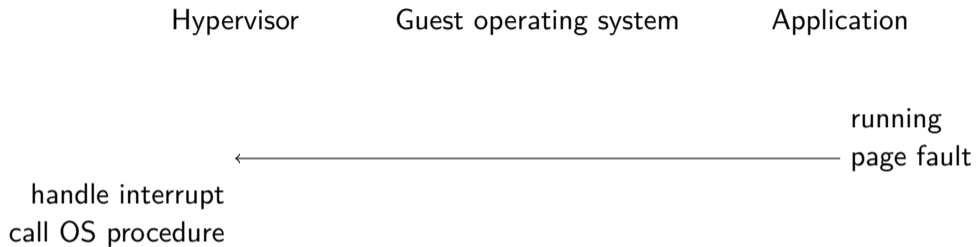
Application



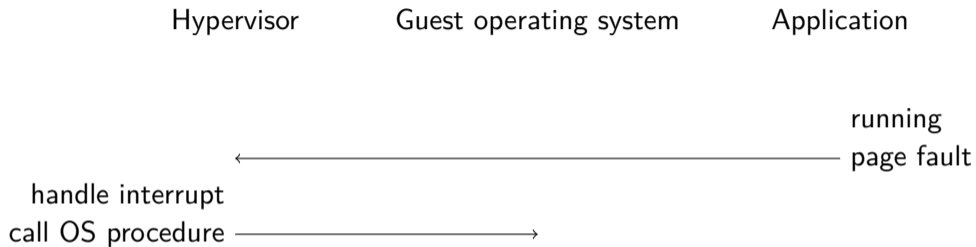
second level paging



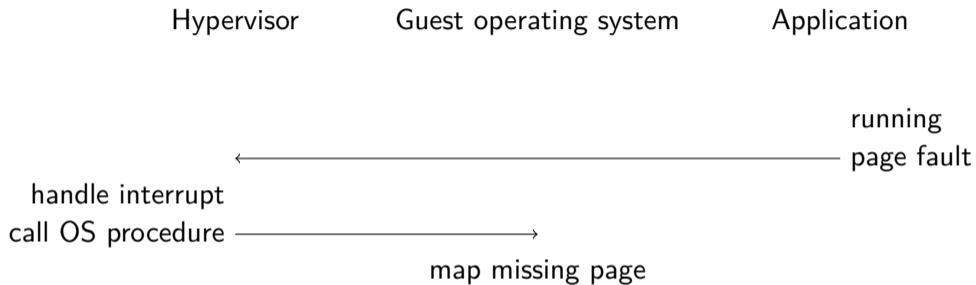
second level paging



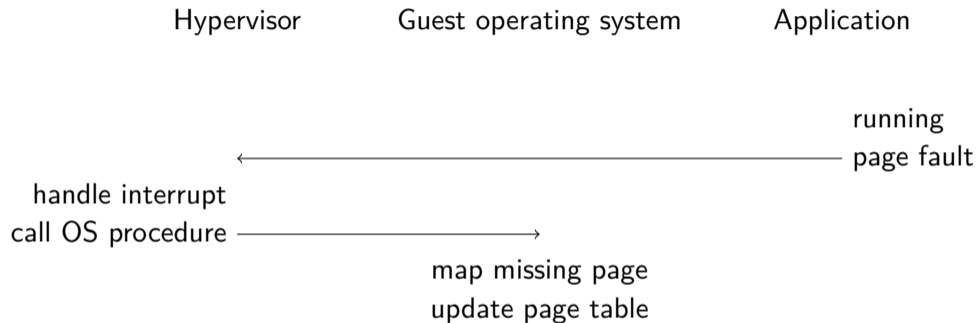
second level paging



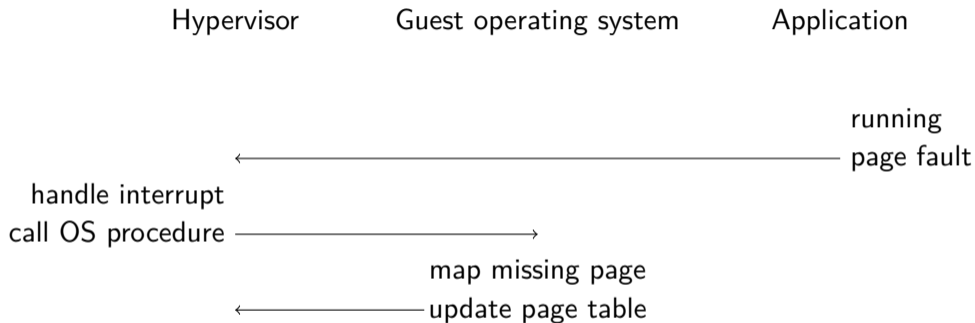
second level paging



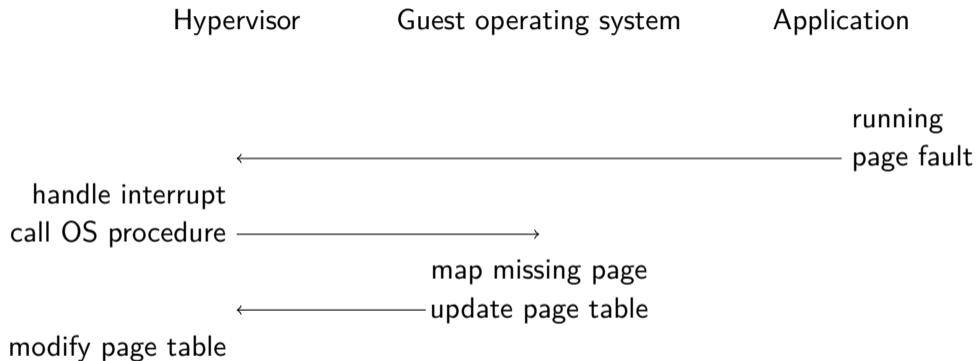
second level paging



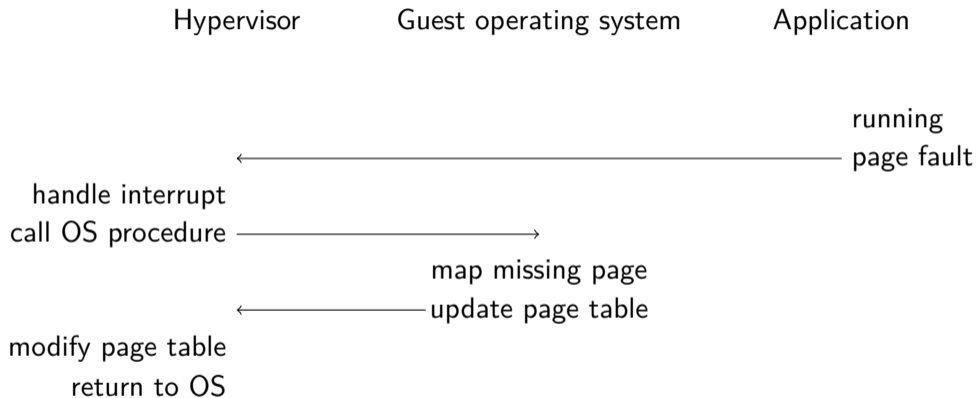
second level paging



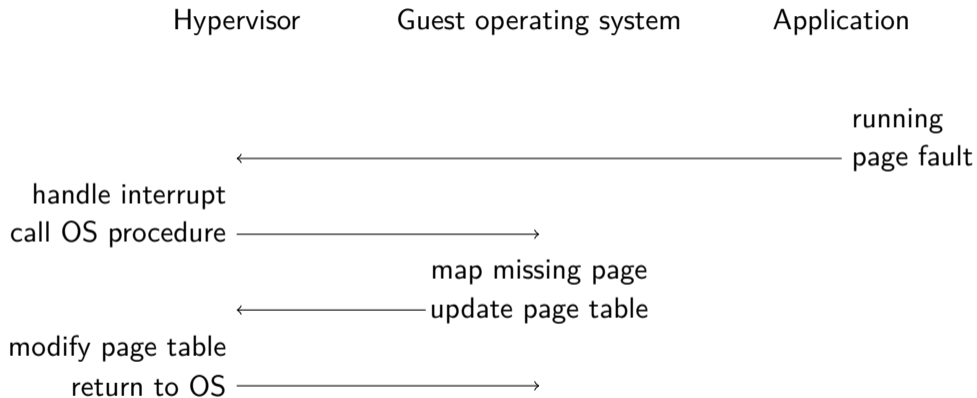
second level paging



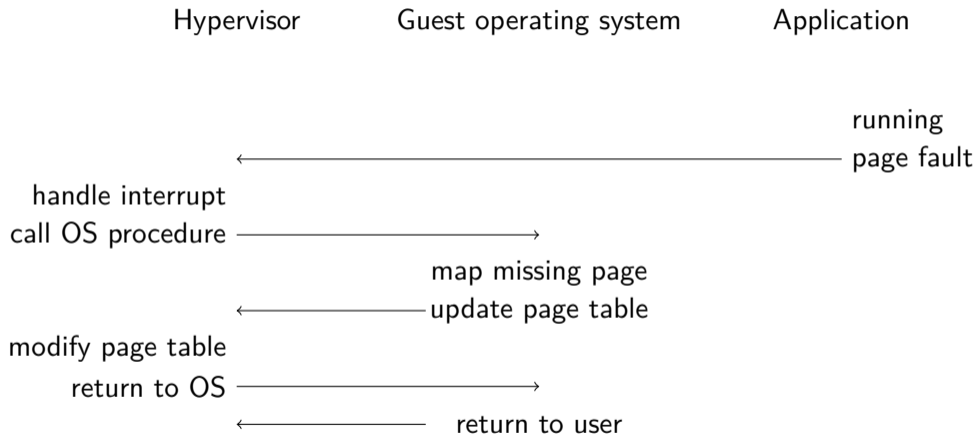
second level paging



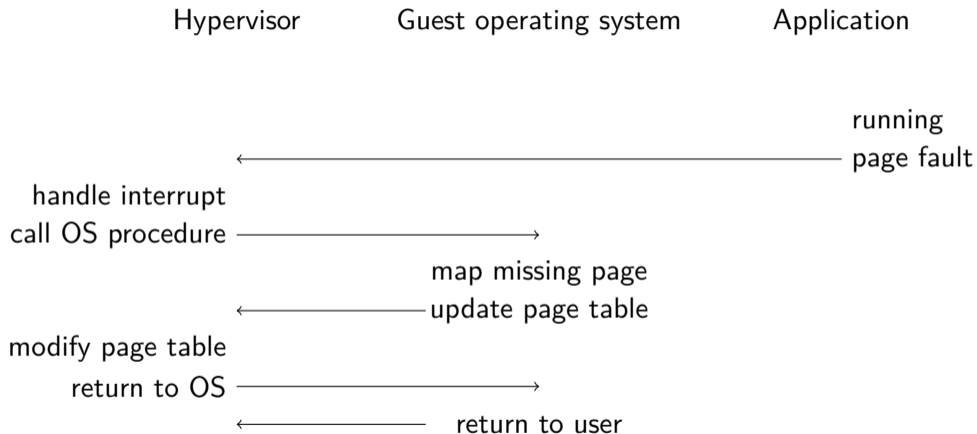
second level paging



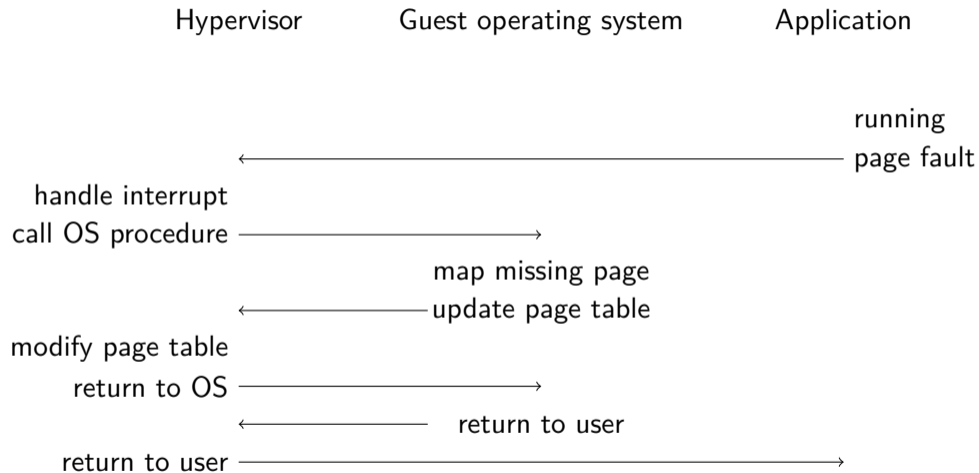
second level paging



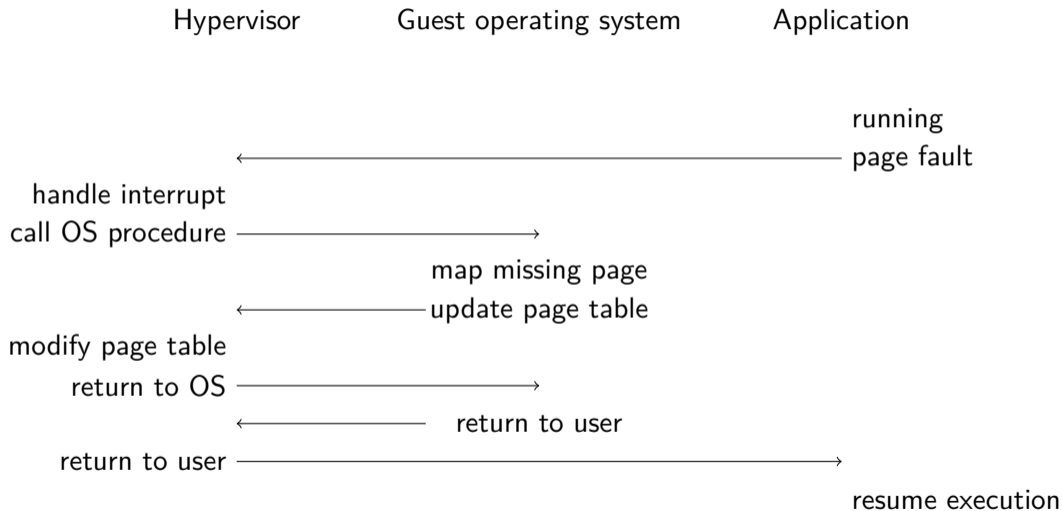
second level paging



second level paging



second level paging



...wait a second

If the guest operating system is executing in user mode - how does it protect itself from the application process that is also running in user mode?

If the guest operating system is executing in user mode - how does it protect itself from the application process that is also running in user mode?

If we allow the guest operating system to run in kernel mode - then the hypervisor can not protect it self.

system call revisited

Hypervisor

Guest operating system

Application

kernel space

user space

in user mode

system call revisited

Hypervisor

Guest operating system

Application

kernel space

user space

in user mode
system call

system call revisited

Hypervisor

Guest operating system

Application

kernel space

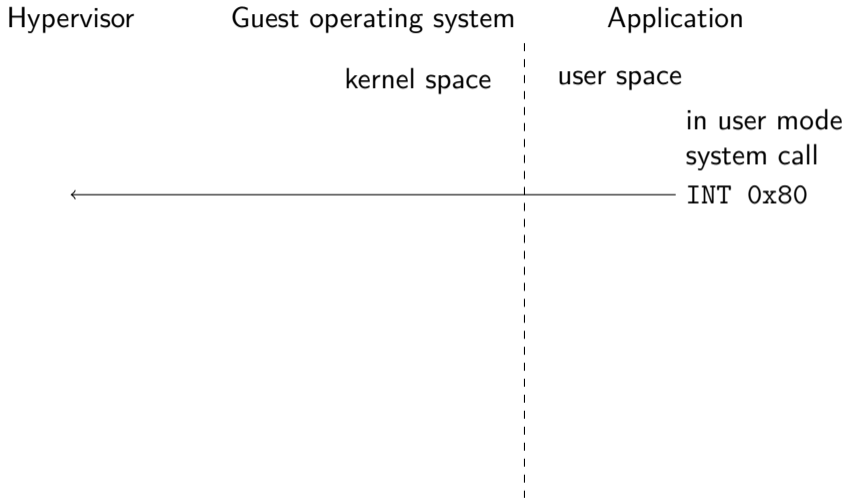
user space

in user mode

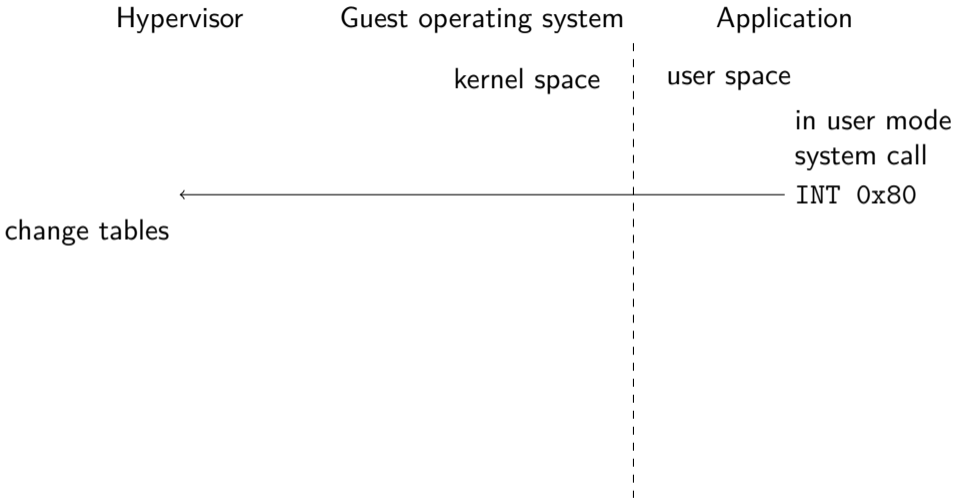
system call

INT 0x80

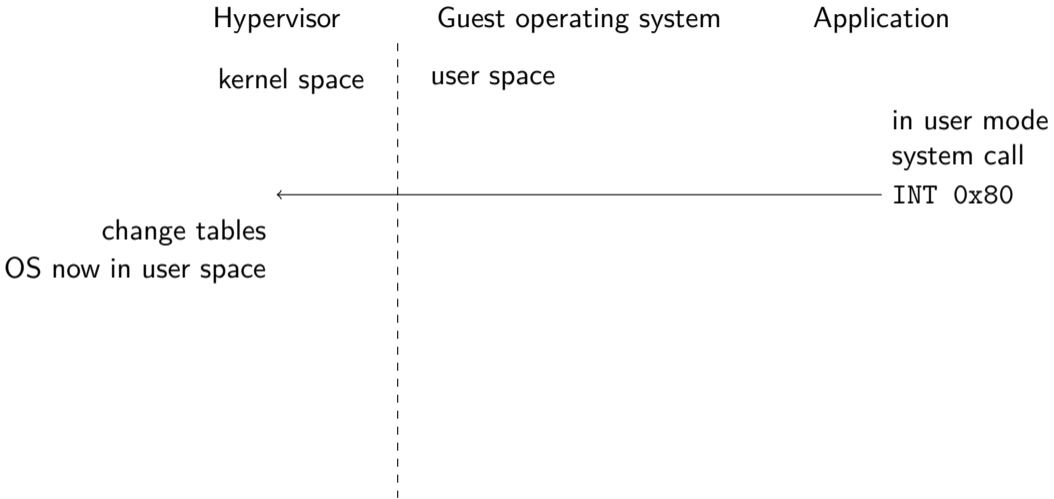
system call revisited



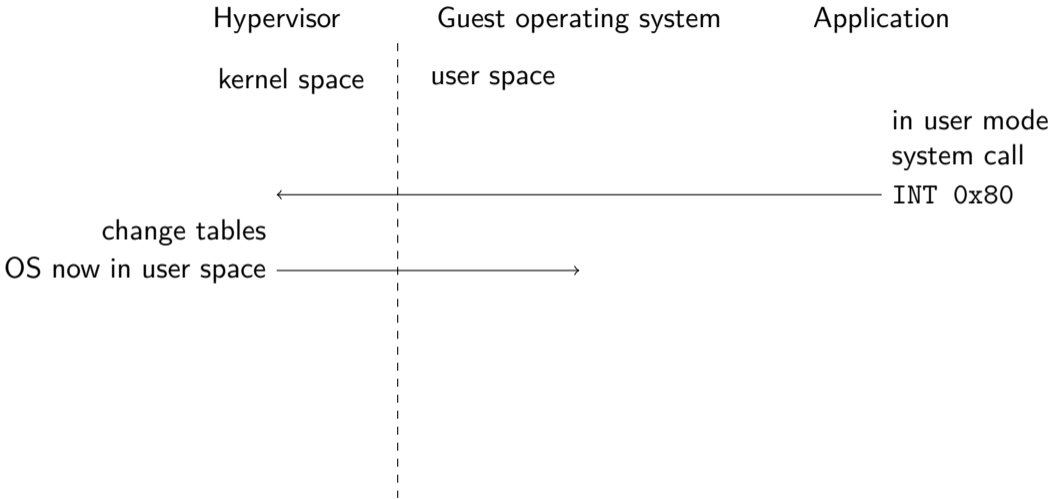
system call revisited



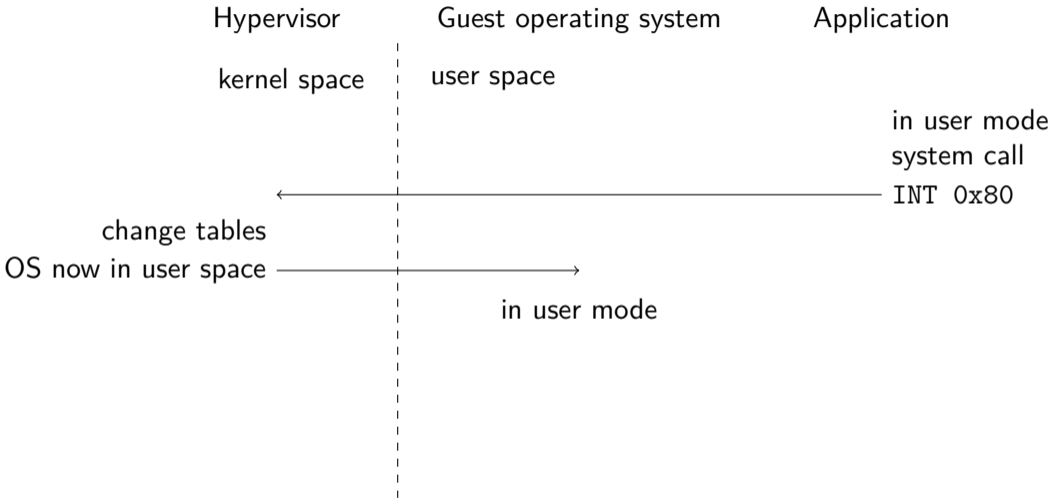
system call revisited



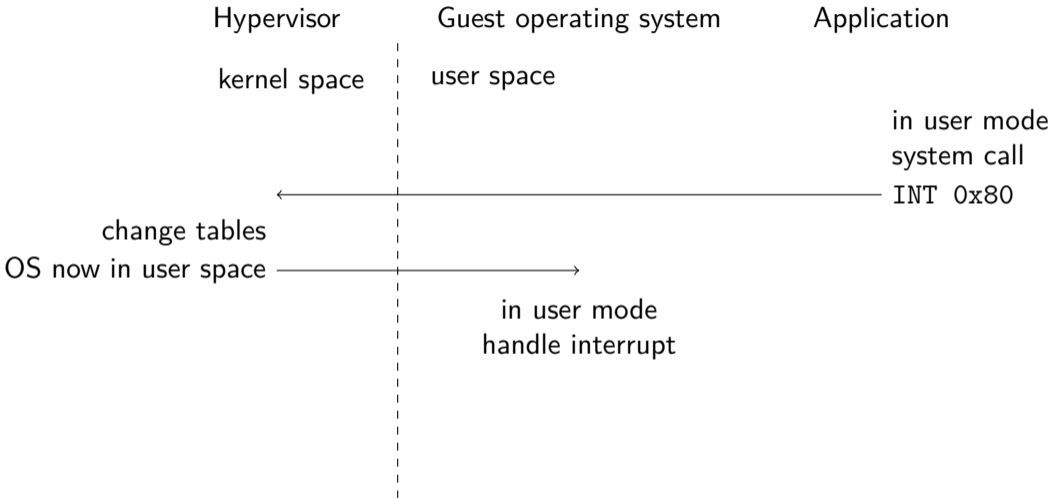
system call revisited



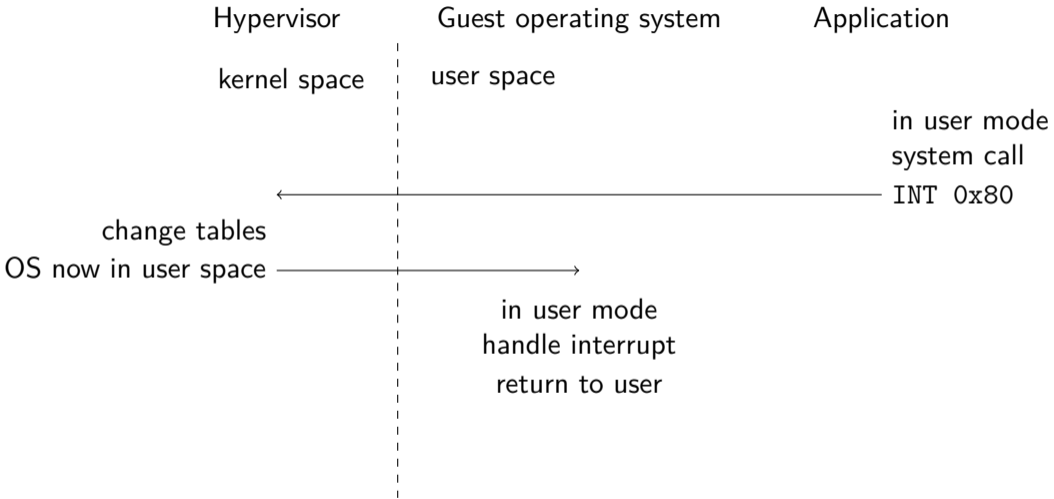
system call revisited



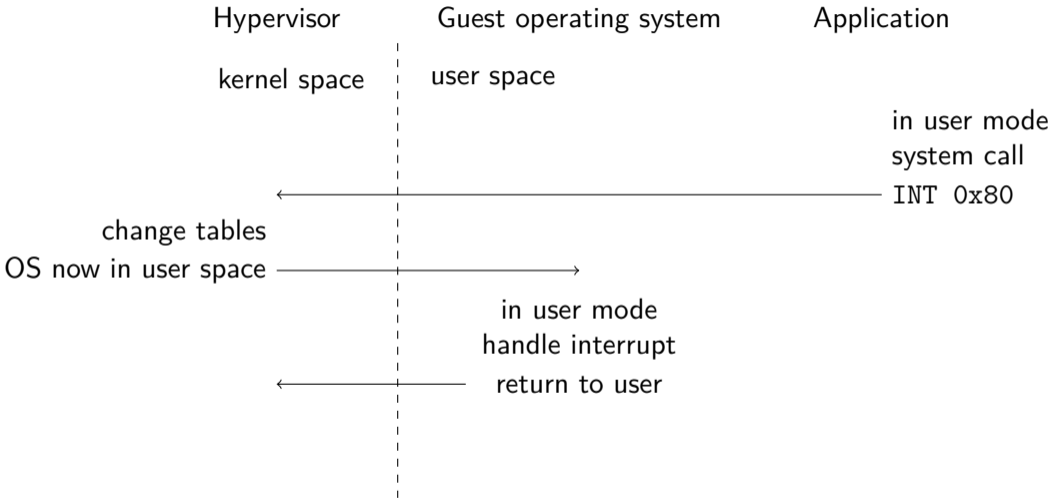
system call revisited



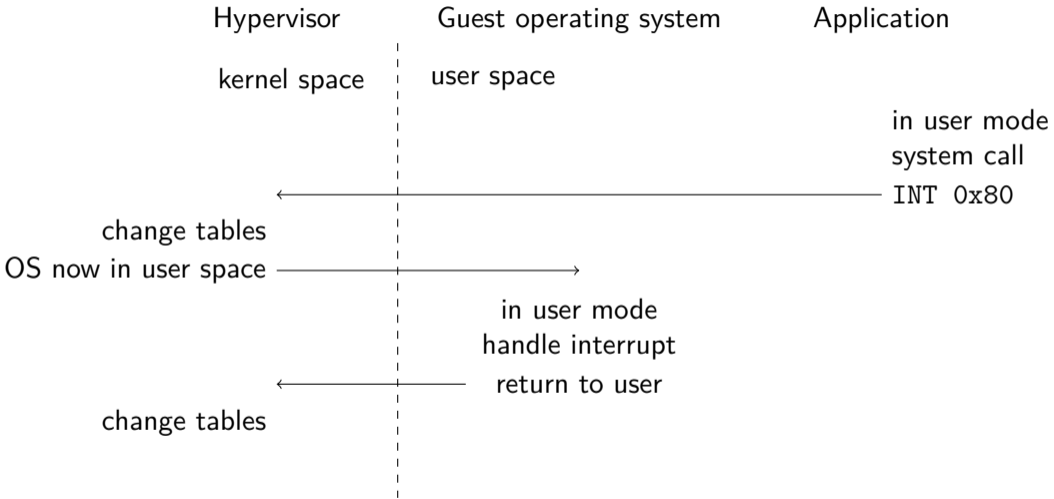
system call revisited



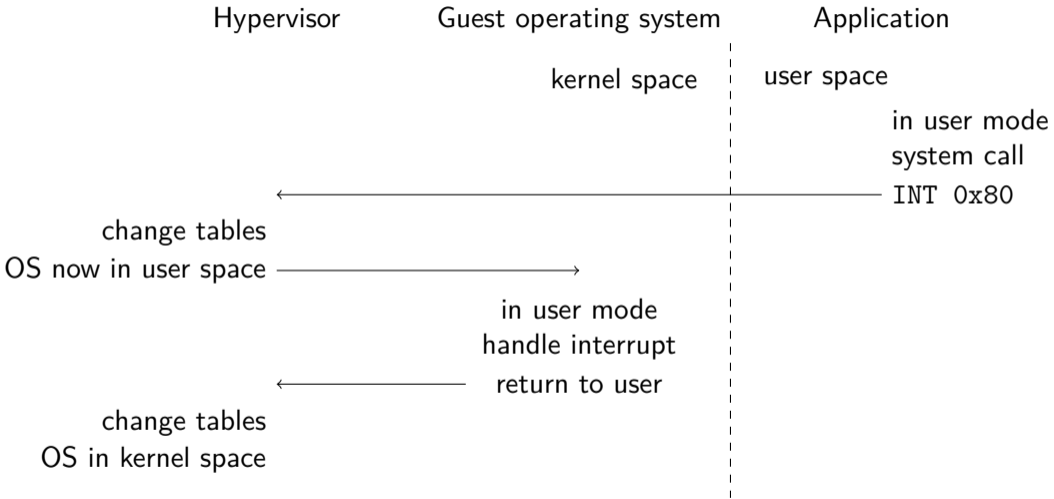
system call revisited



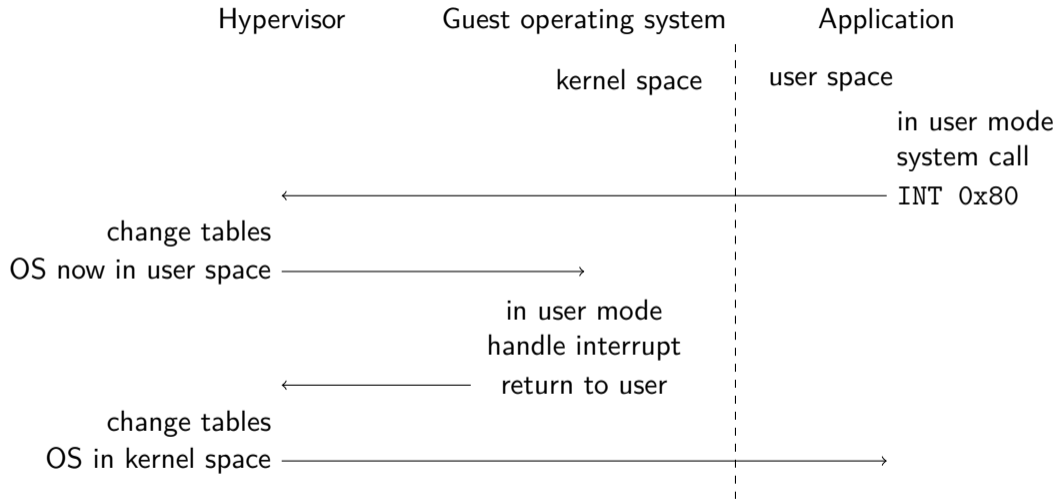
system call revisited



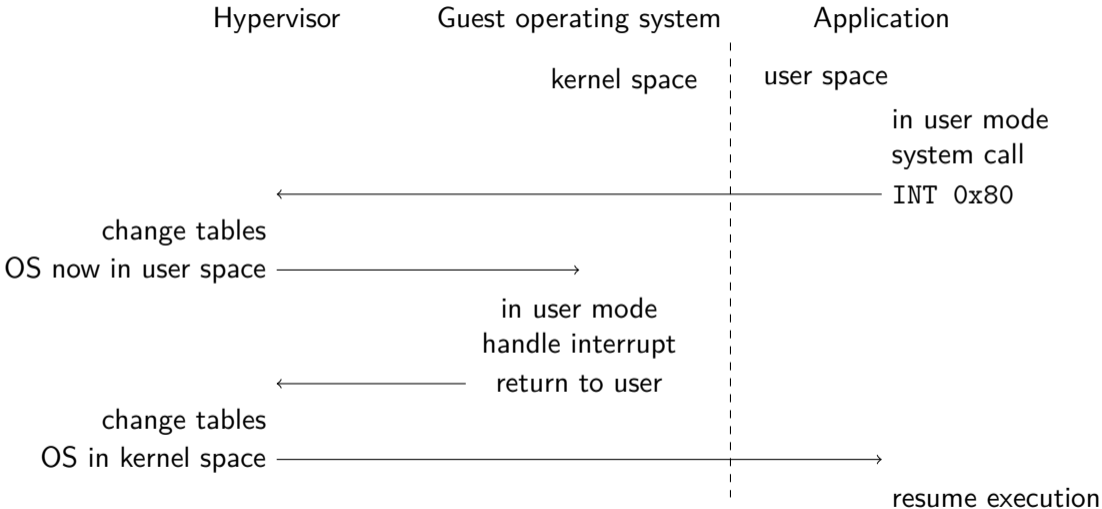
system call revisited



system call revisited



system call revisited



.. thank god for hardware

.. thank god for hardware

Hardware support:

Hardware support:

- Available in both AMD and Intel x86 processors

Hardware support:

- Available in both AMD and Intel x86 processors
- Allows hypervisors to provide near “bare metal” performance.

.. a different approach

Para-virtualization: change the operating system that you want to virtualize.

Para-virtualization: change the operating system that you want to virtualize.

- Change kernel modules in the operating system.

Para-virtualization: change the operating system that you want to virtualize.

- Change kernel modules in the operating system.
- Recompile source code or patch binary code.

The original goal

The original goal

Utilisation of hardware.

The original goal

Utilisation of hardware.

The original goal

Utilisation of hardware.

Applications are completely separated from each other.

The original goal

Utilisation of hardware.

Applications are completely separated from each other.

The original goal

Utilisation of hardware.

Applications are completely separated from each other.

Applications can use different operating systems.

The original goal

Utilisation of hardware.

Applications are completely separated from each other.

Applications can use different operating systems.

What if we skip this.

An operating system uses several name spaces: memory addresses, file paths, port numbers, device interrupt requests, process id, user id, ...

An operating system uses several name spaces: memory addresses, file paths, port numbers, device interrupt requests, process id, user id, ...

Provide a *container*, a separate environment with its own name spaces.

An operating system uses several name spaces: memory addresses, file paths, port numbers, device interrupt requests, process id, user id, ...

Provide a *container*, a separate environment with its own name spaces.

Processes in different containers are completely separated from each other ...

An operating system uses several name spaces: memory addresses, file paths, port numbers, device interrupt requests, process id, user id, ...

Provide a *container*, a separate environment with its own name spaces.

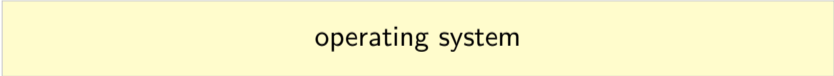
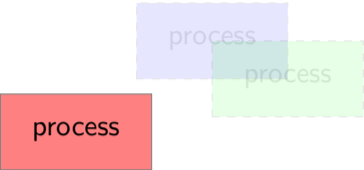
Processes in different containers are completely separated from each other ... but they use the same kernel.

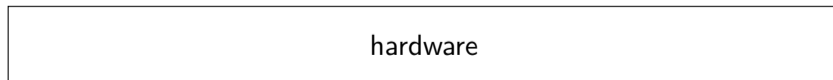
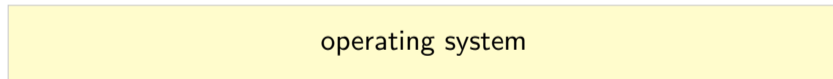
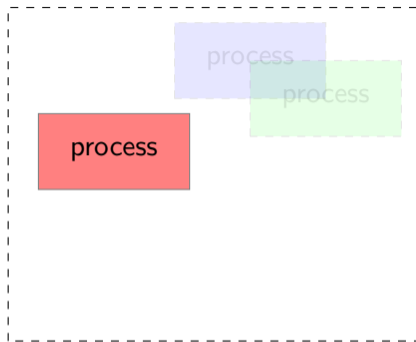
hardware

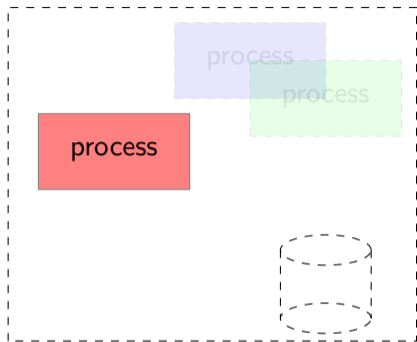
A diagram consisting of two stacked rectangular boxes. The top box is yellow and contains the text 'operating system'. The bottom box is white with a black border and contains the text 'hardware'. The boxes are centered horizontally and stacked vertically.

operating system

hardware



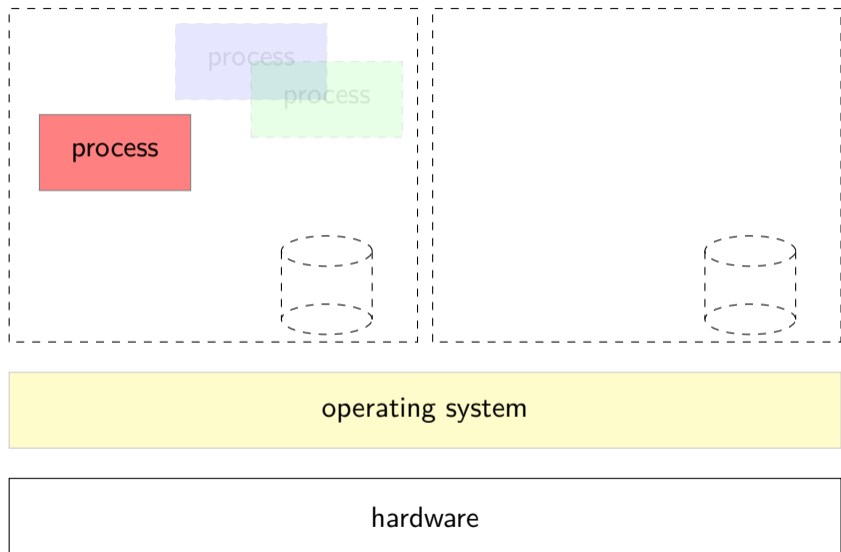




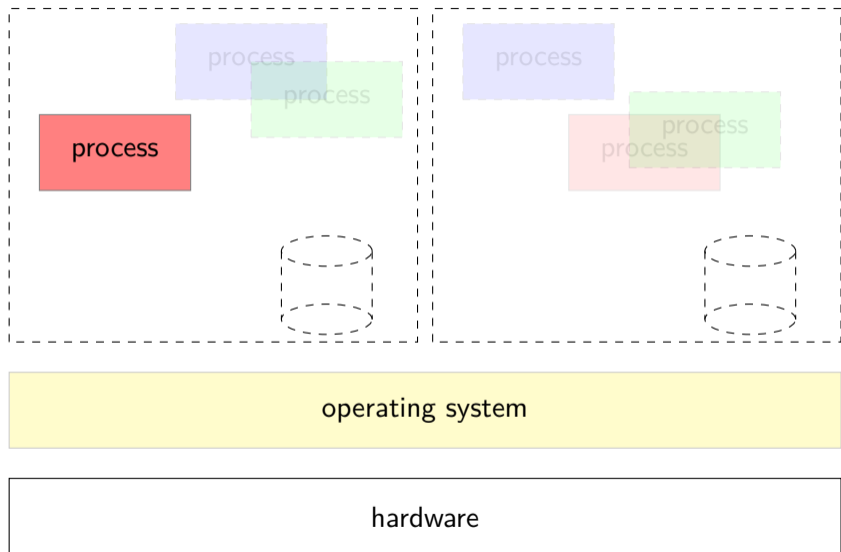
operating system

hardware

containers



containers



the original goal

Utilisation of hardware.

Utilisation of hardware.

Utilisation of hardware.

Applications are completely separated from each other.

Utilisation of hardware.

Applications are completely separated from each other.

Utilisation of hardware.

Applications are completely separated from each other.

Applications can use different operating systems.

Utilisation of hardware.

Applications are completely separated from each other.

Applications can use different operating systems.

Why do they have to run on the same hardware?

x86 hardware

operating system

x86 hardware

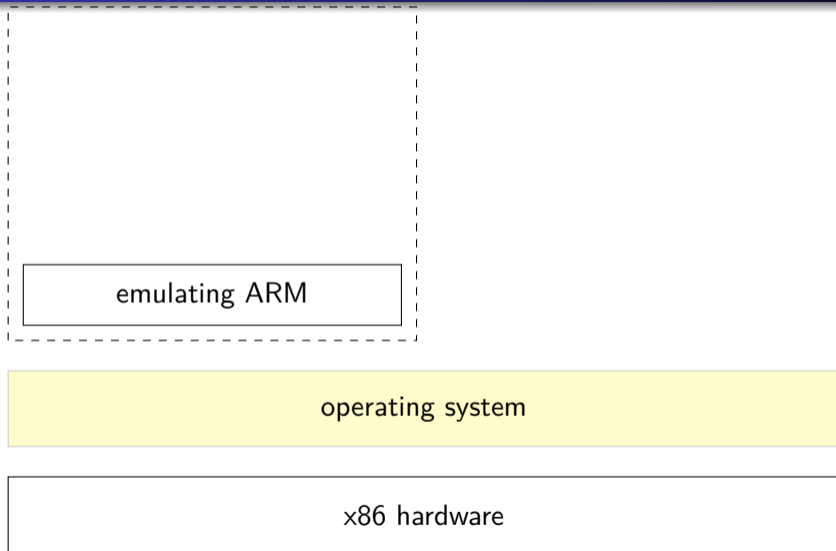
emulating hardware



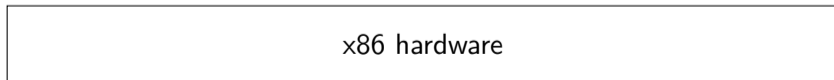
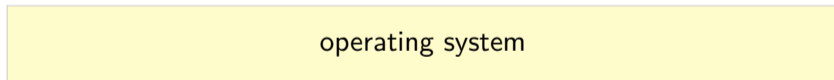
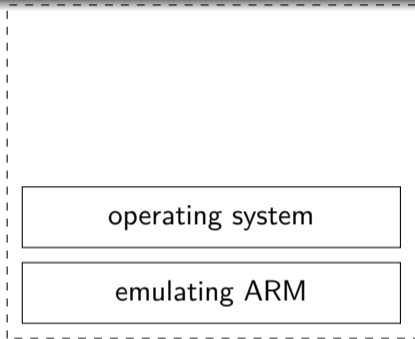
operating system

x86 hardware

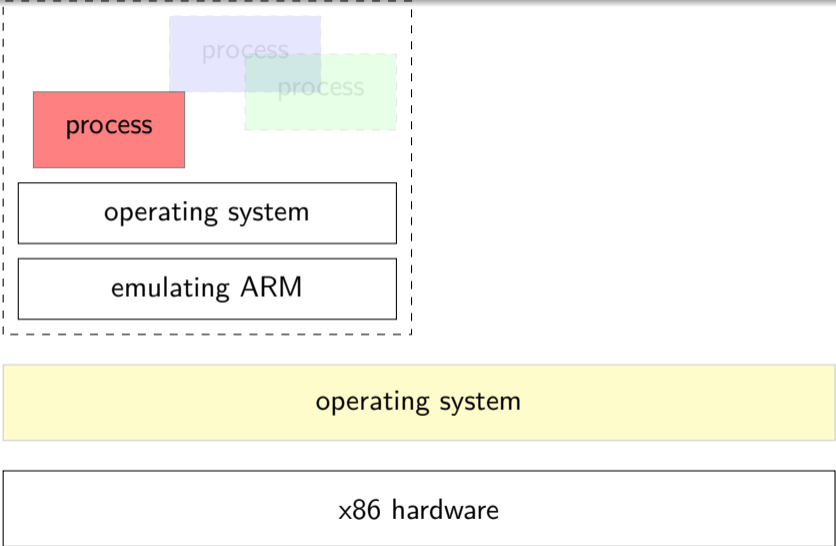
emulating hardware



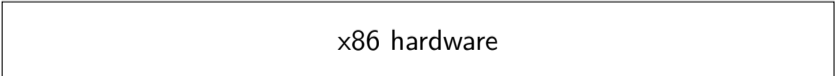
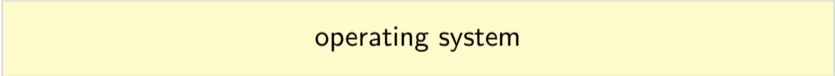
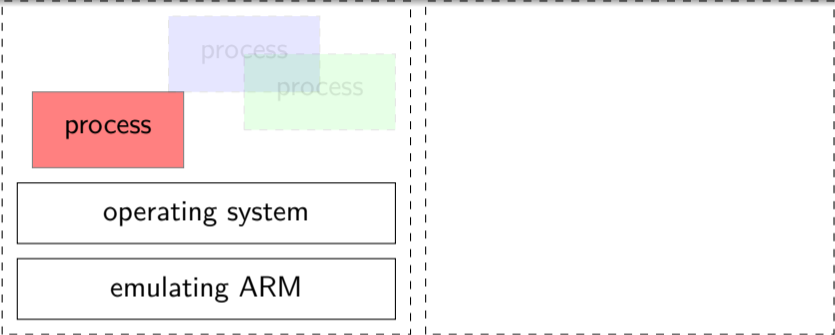
emulating hardware



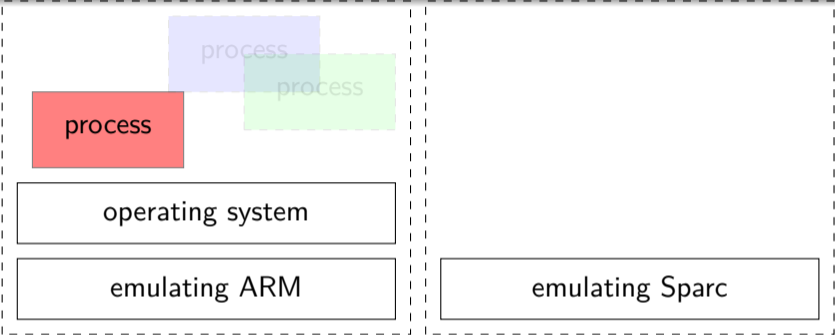
emulating hardware



emulating hardware



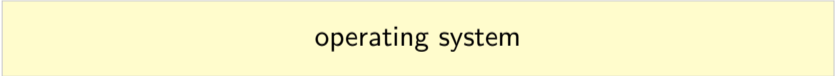
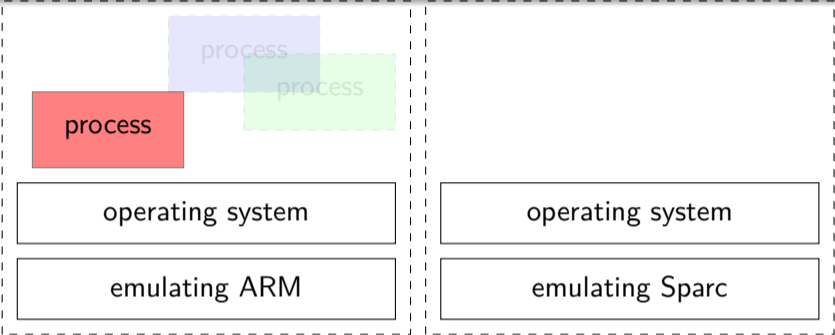
emulating hardware



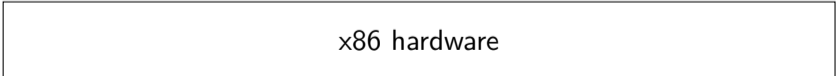
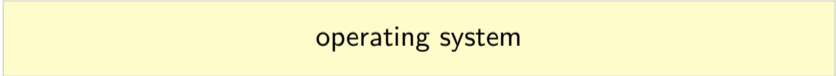
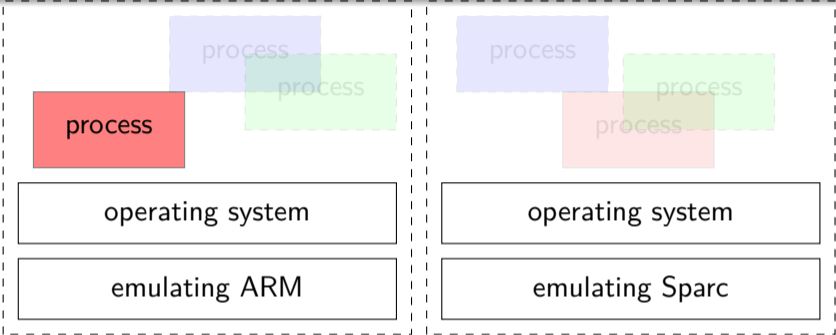
operating system

x86 hardware

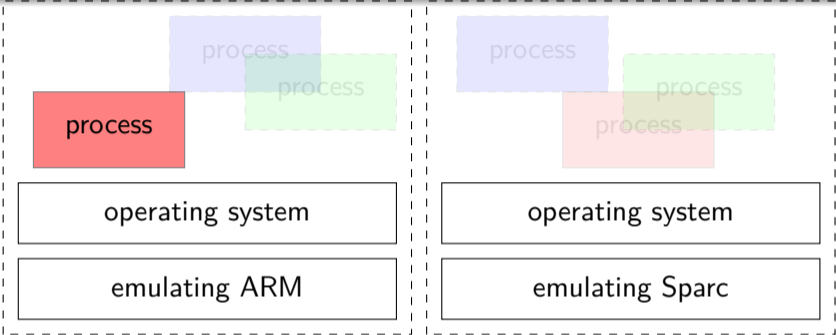
emulating hardware



emulating hardware



emulating hardware



operating system

x86 hardware

- Emulators
 - Can emulate a different hardware than the host machine (QEMU, Simics).

Types of virtual machines

- Emulators
 - Can emulate a different hardware than the host machine (QEMU, Simics).
- Virtual machines
 - Choose operating system but hardware is set (Xen, KVM, VirtualBox, VMware).

Types of virtual machines

- Emulators
 - Can emulate a different hardware than the host machine (QEMU, Simics).
- Virtual machines
 - Choose operating system but hardware is set (Xen, KVM, VirtualBox, VMware).
- Containers
 - Separated name spaces in the same operating system (Docker, Linux Containers).

Types of virtual machines

- Emulators
 - Can emulate a different hardware than the host machine (QEMU, Simics).
- Virtual machines
 - Choose operating system but hardware is set (Xen, KVM, VirtualBox, VMware).
- Containers
 - Separated name spaces in the same operating system (Docker, Linux Containers).
- Runtime systems
 - Dedicated to a language (JVM, Erlang).

... but I never installed a Hypervisor?

... but I never installed a Hypervisor?

VirtualBox etc also installs a kernel module that turns your regular operating system into a hypervisor.

- Multiple operating systems running on the same machine.

- Multiple operating systems running on the same machine.
- Each operating system provided a virtual hardware.

- Multiple operating systems running on the same machine.
- Each operating system provided a virtual hardware.
- With hardware support, near bare metal execution speed can be obtained.

- Multiple operating systems running on the same machine.
- Each operating system provided a virtual hardware.
- With hardware support, near bare metal execution speed can be obtained.
- Other types: emulators, containers, runtime environments.

- Multiple operating systems running on the same machine.
- Each operating system provided a virtual hardware.
- With hardware support, near bare metal execution speed can be obtained.
- Other types: emulators, containers, runtime environments.