

Mechanism: Limited Direct Execution



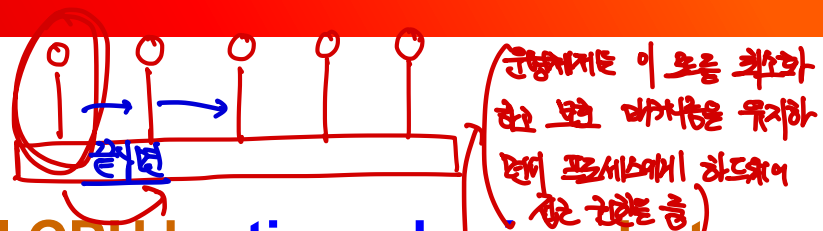
Prof. Yongtae Kim

Computer Science and Engineering
Kyungpook National University

Direct Execution

→ 프로세서 직접 하드웨어를

점진-



- The OS needs to share the physical CPU by **time sharing**, but there are two challenges to virtualize the CPU:
 - Performance: how can we implement virtualization without adding excessive overhead to the system?
 - Control: how can we run processes efficiently while retaining control over the CPU? – control is particularly important to the OS as it's in charge of resources
- The **direct execution** gives rise to a few problems to virtualize
 - How can the OS make sure the program doesn't do anything we don't want?
 - How does the OS stop it from running and switch to another process?



OS	Program
Create entry for <u>process list</u> Allocate memory for program Load program into memory Set up stack with argc/argv Clear registers Execute call main() Free <u>memory</u> of process Remove from <u>process list</u>	Run main() Execute return from main

PCB

"무한프 돌턴...?"

Without **limits** on running on programs, the OS wouldn't be in control of anything and thus would be "**just library**".

Problem #1: Restricted Operations

CPU가 user mode, kernel mode 를 구별해줘야 한다.

What if the process wishes to perform restricted operations?

- Issuing an I/O request to a disk, gaining access to more CPU and memory
- Then, the process can access an entire disk and destroy important data → 파괴하면..?
- Most of operations (arith, loop) can run directly but some should run indirectly

신호

제한은 둘.

운영체제를 통해서.

A new CPU mode, **user mode**, is introduced to restrict access

- Code that runs in user mode is restricted in what it can do (e.g. no I/O request)
- Doing so results in the processor raising an exception and OS kills the process

이 증상

In contrast to user mode is **kernel mode**, which the OS runs in

- Code that runs can do what it likes including privileged operations (e.g. I/O)

Thus, the CPU must support at least two mode of operations

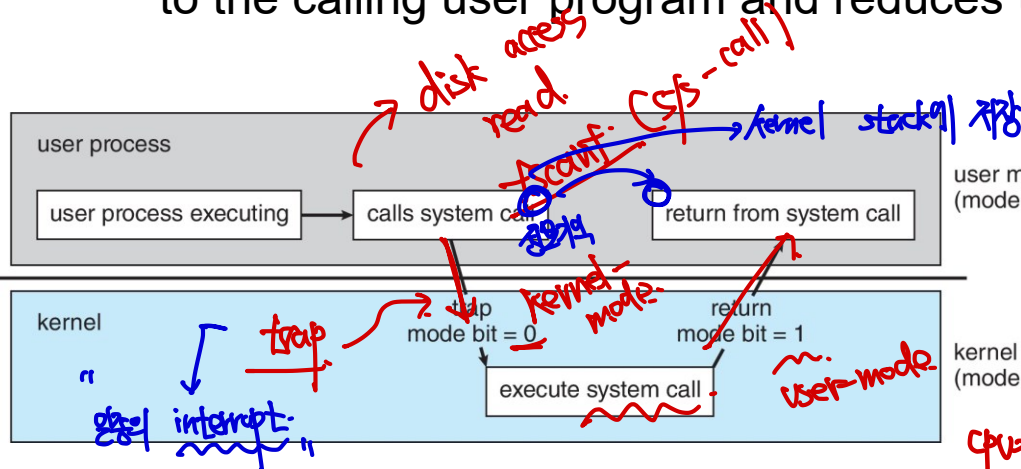
더 많은 것

- e.g.) 4 privilege levels in x86 and its level is set by current privilege level (CPL) in CS register
- The privileged instructions can only be executed in the corresponding level; otherwise, the CPU raises an exception and the OS kills the process

System Call and Trap

fork() → 실행 트랩 발생하기 CPU 실행

- What if a **user program** wants to perform **privileged operations**?
 - For this, **hardware** provides the ability for user programs to perform a **system call**.
 - System calls allow the kernel to **carefully expose certain key pieces of functionality** to user programs, such as accessing file system, allocating more memory
- For **system call**, a program must execute a special **trap** instruction
 - This simultaneously **jumps into the kernel** and raises the **privilege level to kernel mode**; Then, the system can perform the privilege operations
 - When finished, the OS calls a special **return-from-trap** instruction that returns to the calling user program and reduces the privilege level to user mode



user mode (mode bit = 1)

- Trap is known as **software interrupt**

kernel mode (mode bit = 0)

- In x86, **"INT"** is trap instruction and **"IRET"** is return-from-trap instruction

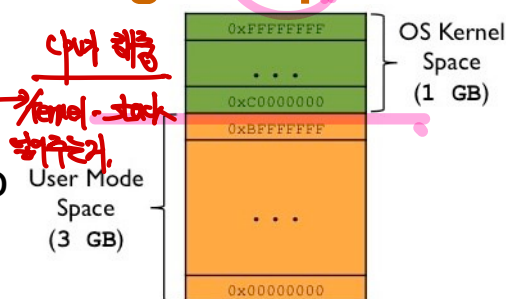
CPU 마다 다른 resource 접근

하드웨어 인터럽트 ↔ hardware interrupt

Trap Handling

The hardware needs to be careful to when executing a trap

- It must save caller's registers to be able to return correctly when the OS issues the return-from-trap instruction
- In x86, the CPU pushes PC, flags, and some registers onto a per-process kernel stack; the return-from-trap pops these values and resumes user-mode program



How does the trap know which code to run in OS?

- The kernel does so by setting up a trap table (at boot time)
- When the machine boots up, it's in privileged (kernel) mode and the OS initializes the trap table
- The OS informs the CPU of the locations of these trap handlers

process user space

이제

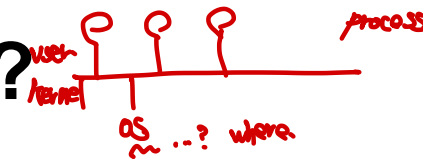
하단 테이블에 해당하는 trap handler 번호를 취하는 것

INT #	Description
0x00	Divide by 0
0x01	Reserved
0x02	NMI Interrupt
0x03	Breakpoint (INT3)
0x04	Overflow (INT0)
0x05	Bounds range exceeded (BOUND)
0x06	Invalid opcode (UD2)
0x07	Device not available (WAIT/FWAIT)
0x08	Double fault
0x09	Coprocessor segment overrun
0x0A	Invalid TSS
0x0B	Segment not present
0x0C	Stack-segment fault
0x0D	General protection fault
0x0E	Page fault
0x0F	Reserved
0x10	x87 FPU error
0x11	Alignment check
0x12	Machine check
0x13	SIMD Floating-Point Exception
0x14-0x1F	Reserved
0x20-0xFF	User definable

A system-call number is usually assigned to each system call

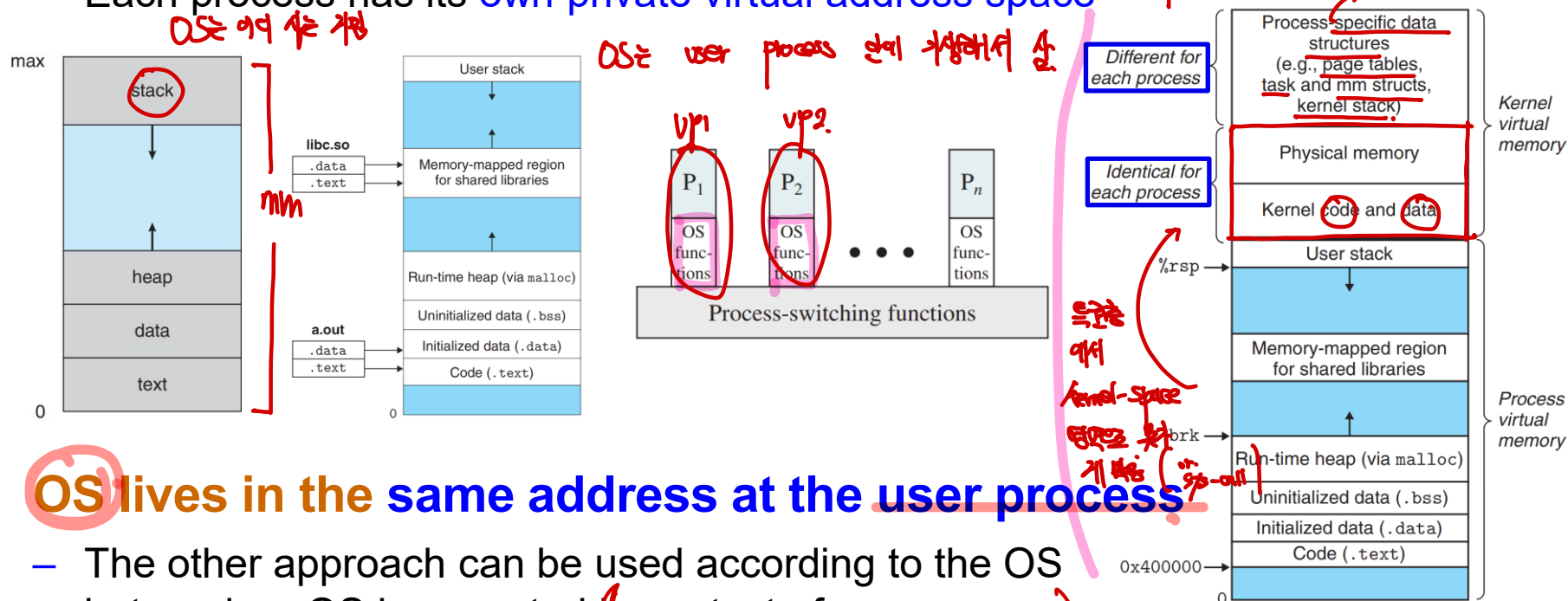
- The OS, when handling the system call inside the trap handler, examines this number, ensures it is valid and executes the codes; a form of protection
- Informing the CPU of trap table location (IDTR) is done by a privileged operation

Where does the OS live?



- **Memory Layout Revisited:** The memory layout of a process is typically divided into four segments of text, data, heap, and stack

- Each process has its own private virtual address space

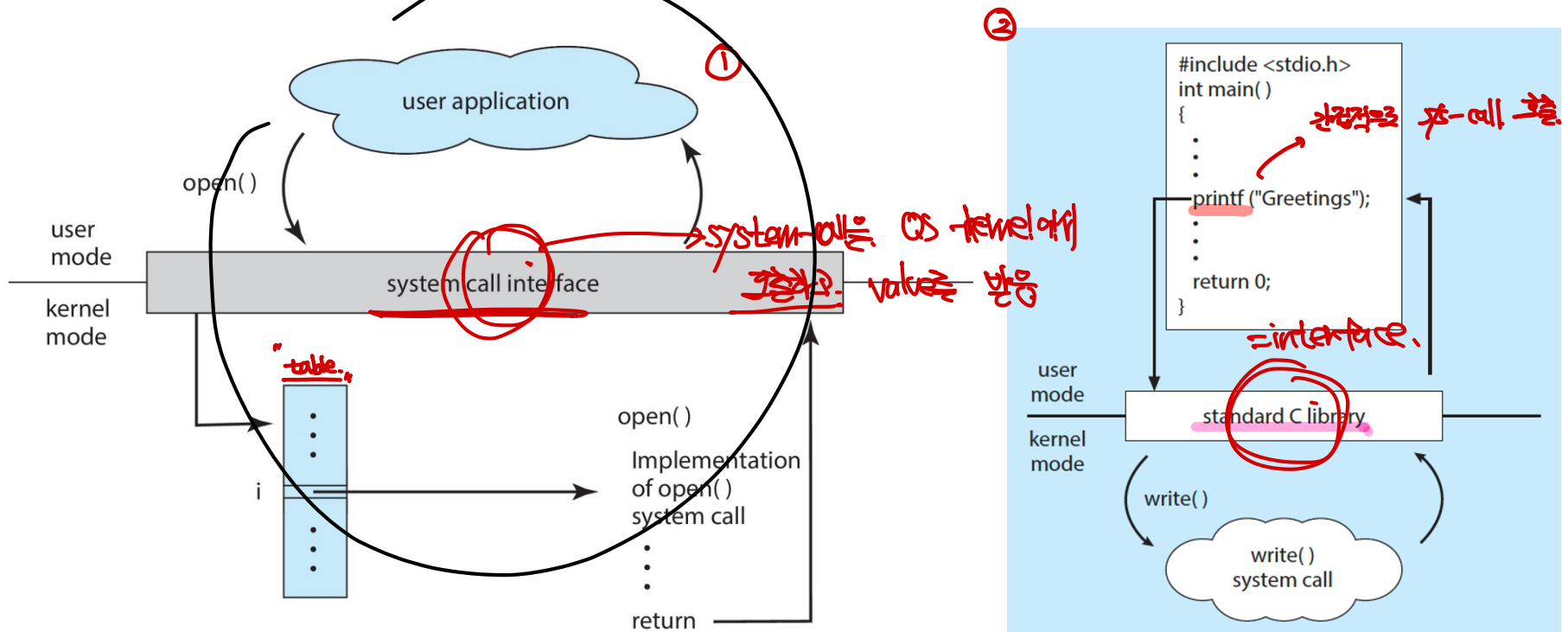


- **OS lives in the same address at the user process**

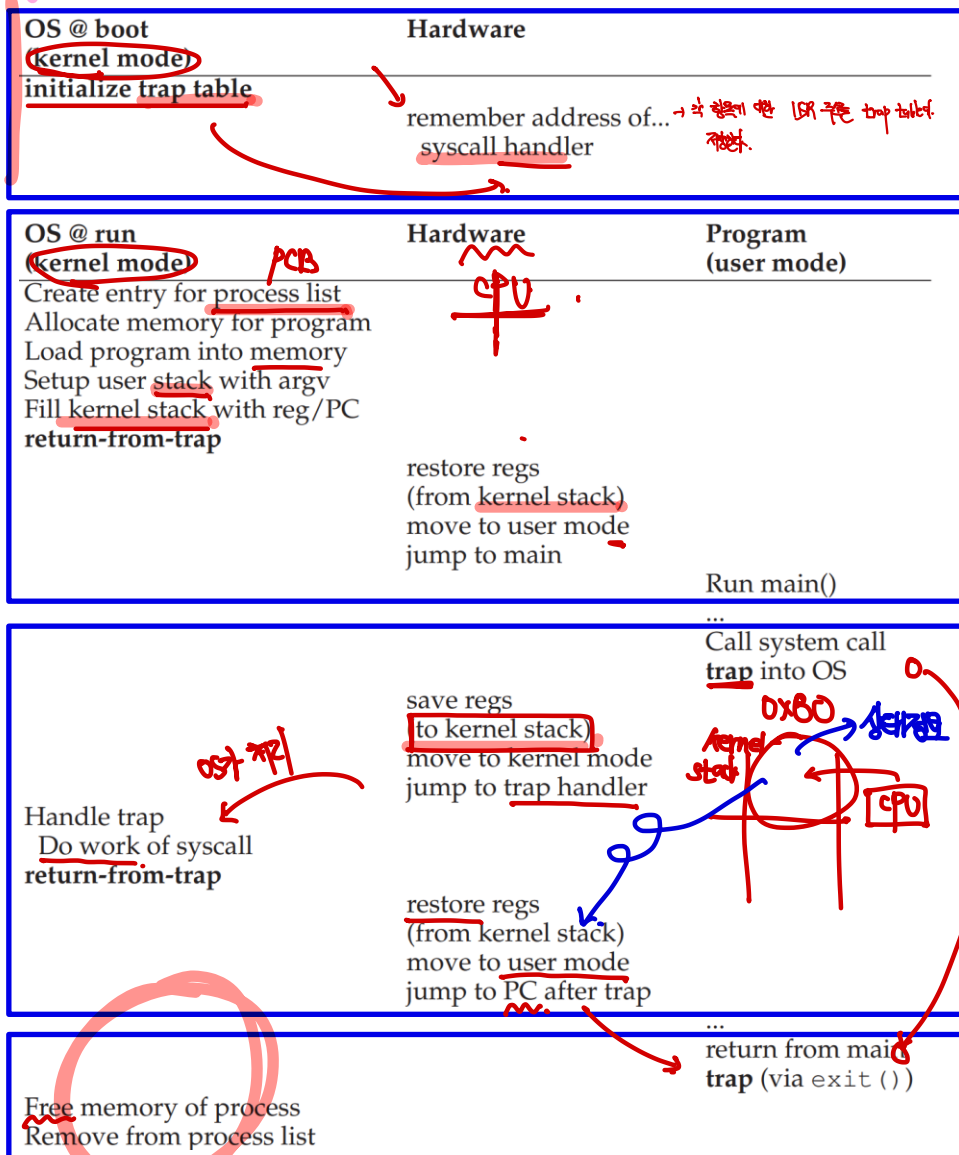
- The other approach can be used according to the OS but modern OS is executed in context of user process.
- The privileged mode does not allow the user process to jump into the kernel space and prevent it from accessing the kernel space
- User programs must access kernel code and data indirectly via the system call

System Call Interface and Standard Library

- **System call interface maintains a table indexed to the numbers**
 - The system call interface invokes the intended system call in OS kernel and returns status of the system call and any return values
- **The standard library provides a portion of system call interface**
 - e.g.) invoking `printf()` ^{wrapping} → intercepting this call by lib → invoking syscall `write()` → taking the return value from `write()` → passing it back to user program



Limited Direct Execution Protocol ^{LDE}



Initializing when booting up

Creating process

Processing system call by trap

- system call → trap → save context and switch stack → jump to the trap handler → processing in kernel mode
- return-from-trap → switch stack and restore context → jump to the next of the system call → continue in user mode

Destroying process

OS
kernel-mode

HW

program.

initialize
trap-table → remember IDTR
trap table
address.

Create entry process
list:

↑
Allocate "memory" to p.

Load prog. into mem.

Set user-stack with arg.

Fill kernel stack with
pc/regs
etc. stuff

return-from-trap

→ pop kernel-stack
restore regs
move to user-mode
jump to mem

→ Run main()

Call sys-call

→ trap into OS

Save (regs) into
kernel-stack

kernel-mode

jump to trap handler

handle trap

do work for sys-call

return-from-trap



restore reg
user code

jump to code after trap



return from main

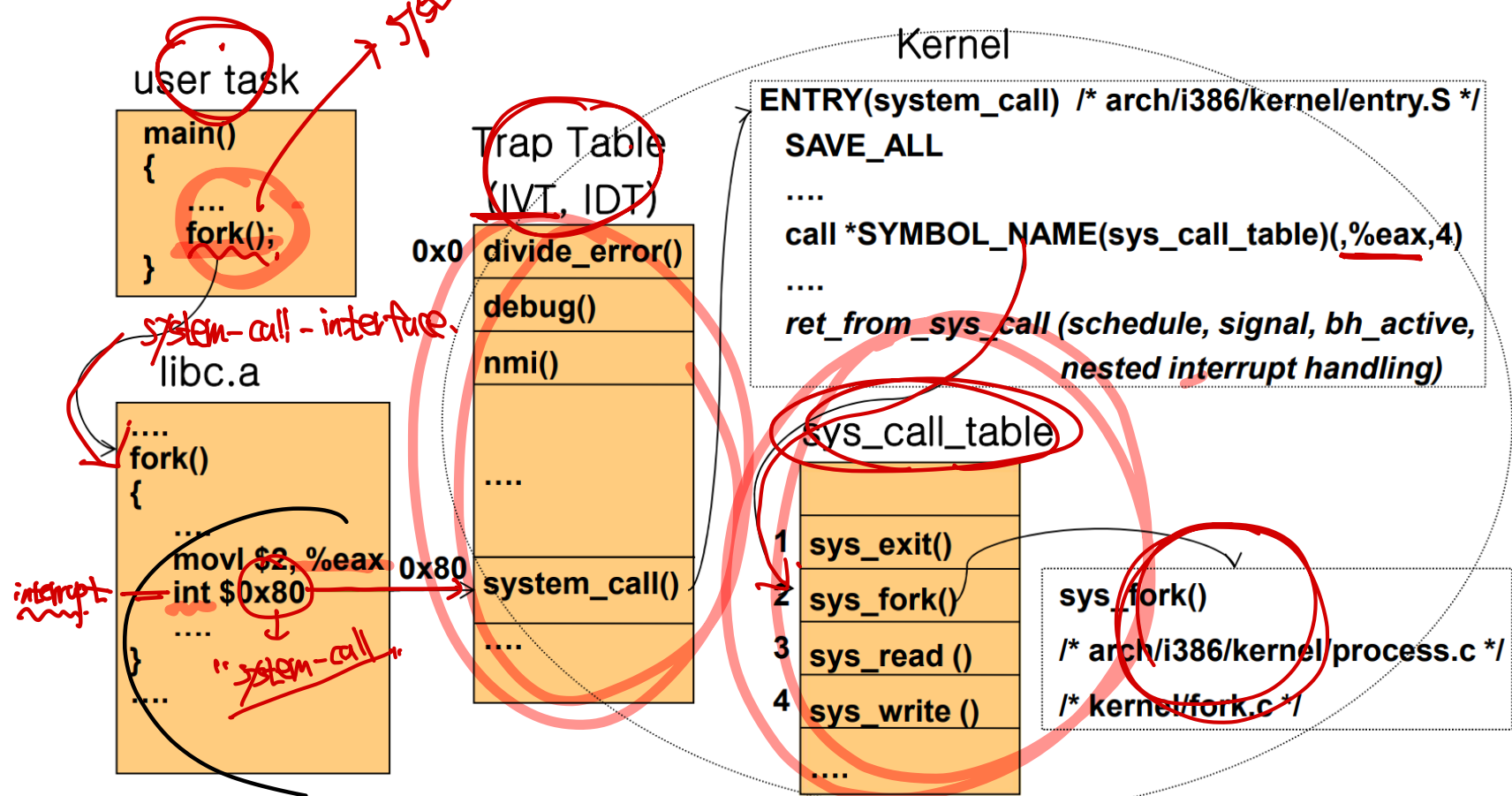
trap



lose

System Call Implementation

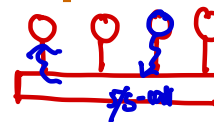
Linux system call on x86



Courtesy of Prof. Jongmoo Choi @ Dankook Univ.

Problem #2: Switching between Processes

- If a process is running on the CPU, then, the OS is not running
 - Then, there is no way for the OS to do for switching between processes
 - How can the OS regain control of the CPU so it can switch between processes?
- One way is a cooperative approach: wait for system calls → 안됨
 - Processes are assumed to periodically give up the CPU so the OS can decide to run some other task
 - Most processes transfer control of the CPU to the OS frequently by system calls
 - Application also transfer control to the OS when they do something illegal, such as dividing by zero, illegal memory access, which will cause a trap
 - yield system call is provided in case processes seldomly use a system call
 - Big issue: what if a process gets stuck in an infinite loop? → reboot the machine
- In short, a cooperative approach makes three steps
 - 1) processes use a system call
 - 2) control transfer to the OS
 - 3) do scheduling (switching)



A Non-Cooperative Approach: The OS Takes Control

- **How can the OS gain control of the CPU without cooperation?**

Hardware interrupt
 timer device + Hardware all exist

 - The answer is simple and was discovered long time ago: a timer interrupt
 - A timer device can be programmed to raise an interrupt every so many millisec.
 - When 1) the interrupt is raised, 2) the running process is halted and 3) a pre-configured interrupt handler in the OS runs → the OS has regained control

- **The OS must (inform the hardware) of which code to run when the timer interrupt occurs**
 - At boot time, the OS do that and must start the timer by a privileged operation

- **The hardware has some responsibility when an interrupt occurs**
 - In particular to save enough of the state of the running process such that a subsequent return-from-trap instruction will be able to resume it correctly
 - This is similar to the behavior of the hardware during an explicit system-call trap into the kernel

Saving and Restoring Context

- **The scheduler, a part of OS, makes a decision whether**

- To continue running the currently-running process OR switch to a different one
- If the decision is made to switch, the OS execute a low-level piece of code, which is referred to as a context switch.

- **A context switch is conceptually simple**

- A few register values for current process are saved onto the memory and a few for the soon-to-be-executing process are restored from the memory
- By doing so, the OS ensures that when return-from-trap instruction is executed, the system resumes execution of another process, instead of the original one

- **The OS will execute some low-level assembly code to**

- save general purpose registers, PC, and the kernel stack pointer of currently-running process
- then, restore the registers and PC, and switch to the kernel stack (by changing kernel stack pointer) for the soon-to-be-executing process



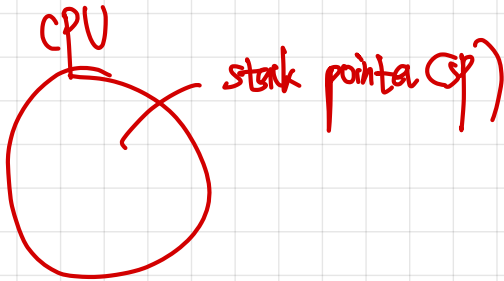
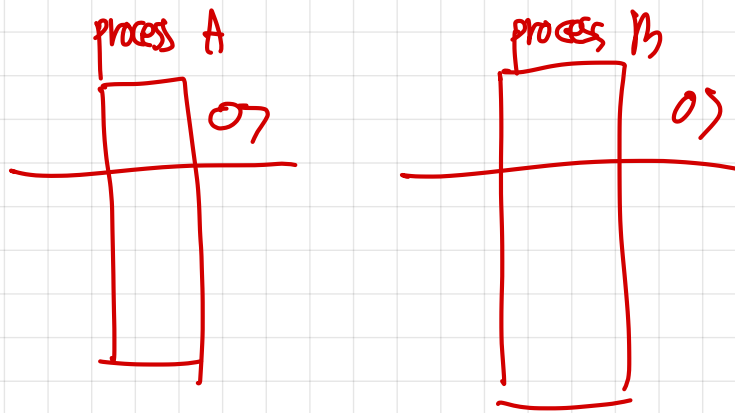
Limited Direct Execution Protocol (Timer Interrupt)

OS @ boot (kernel mode)	Hardware
initialize <u>trap table</u>	remember addresses of... <u>syscall handler</u> <u>timer handler</u>
start <u>interrupt timer</u>	start <u>timer</u> interrupt CPU in X ms

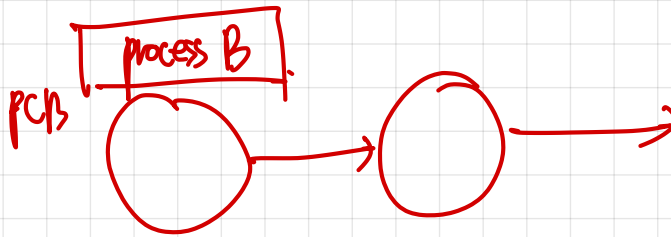
OS @ run (kernel mode)	Hardware	Program (user mode)
		Process A ...
	timer interrupt save regs(A) → k-stack(A) move to kernel mode jump to trap handler	
Handle the trap Call switch() routine save regs(A) → proc.t(A) restore regs(B) ← proc.t(B) switch to k-stack(B) return-from-trap (into B)		
	restore regs(B) ← k-stack(B) move to user mode jump to B's PC	Process B ...

- Initializing when booting up**
- Running process A & Timer interrupt**
- Scheduling & Context switch (A → B)**
 By changing the stack pointer to use the B's kernel stack (not A's)
- Returning from interrupt & Running process B**

- There are two types of register saves and restores in LDE**
 - The user registers of the running process are implicitly saved to its kernel stack by hardware (CPU) when a timer interrupt occurs
 - The kernel registers are explicitly saved to its PCB in memory by software (OS) when the OS decides to switch



kernel



IPT IRT
trap table

1. timer interrupt
2. (CPU) regs(A) → kernel (A)
↓ (kernel-mode)
- 3.

OS boot.
(kernel-mode)

Hardware.

initialize trap table.

(remember address of
syscall-handler,
timer-interrupt
handler)

start
interrupt timer

(Xmsec 지정 CPU를 인터럽트)
↓
호출

운영체제
(커널모드)

호출

프로세스

↓
프로세스 A
...

OS에 요청으로
타이머 인터럽트
↓

A의 레지스터를 kernel-stack에 저장

가장 먼저

연속 호출

연속 호출

switch C를

A의 요청에 대해 A의 proc에 저장

A의 proc에 저장된 A의 레지스터 복원

A의 kernel-stack을 복원

return-from-trap

A의 가장 먼저 A의 레지
스택에 저장

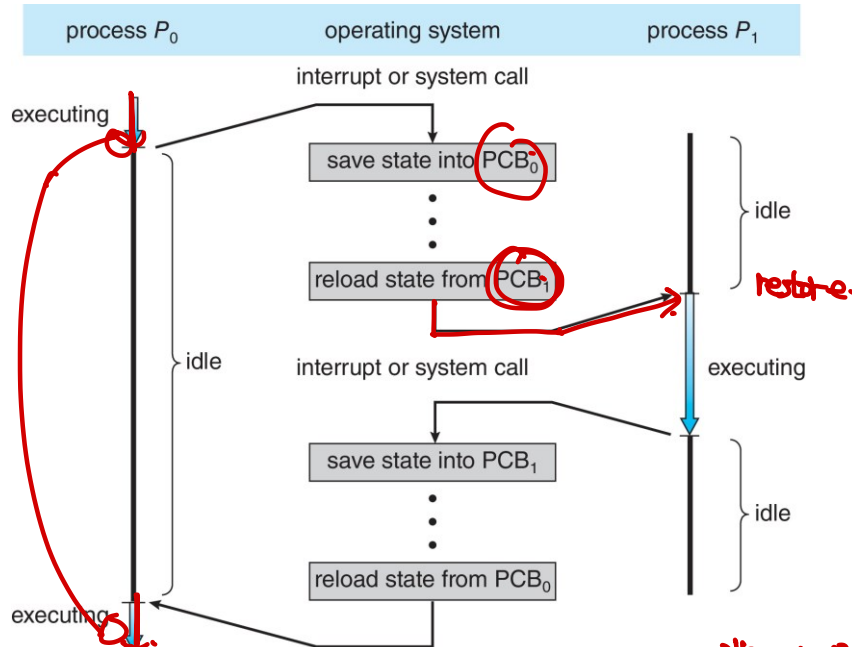
시작

A의 proc

Context Switching Revisited and Concurrency Issue

CPU가 쉬고 있을 때만 가능.

- The context of a process is represented in the PCB



- Statistics of context switches in Linux

→ 커널 값 반환

- Total uptime: 2,176,310 sec (25 days)
- Total 14,115,783,018 context switches
- Average 6486 context switches/sec
- Average 541 context switches/sec/CPU (for all 12 CPUs)

12 core

```
[tahoe: /home/yongtae]$ cat /proc/uptime
2176310.83 66103492.27
[tahoe: /home/yongtae]$ grep ctxt /proc/stat
ctxt 14115783018
[tahoe: /home/yongtae]$ echo "14115783018 / 2176310.83" | bc -l
6486.10613126434701425439
[tahoe: /home/yongtae]$ echo "6486.10613126434701425439 / 12" | bc -l
540.50884427202891785453
[tahoe: /home/yongtae]$
```

idle 한 상태가 많았을 때

HW-interrupt

- Worried about concurrency?

- What happens when, during a system call, a timer interrupt occurs?
- What happens when you're handling one interrupt and another one happened?
- Then, the OS can disable interrupt during the interrupt handling (dangerous), and use some sophisticated locking schemes → concurrency issue