# *ECE391*
# *Computer System Engineering*
### *Lecture 19*

Dr. Jian Huang

University of Illinois at Urbana- Champaign

Fall 2018

# *ECE391 EXAM 1*

- **EXAM II – November 6; 7:00pm-9:00pm**
  - Seating:
    - ECEB 1002: last name starting with A to W
    - ECEB 1013: last name starting with X to Z

- **Review Session**
  - In Class; Thursday November 1

- **Conflict Exam**
  - Tuesday November 6, 2:00pm; Location: CSL 141
  - By Friday November 2 send email to kalbarcz@Illinois.edu to request Conflict Exam
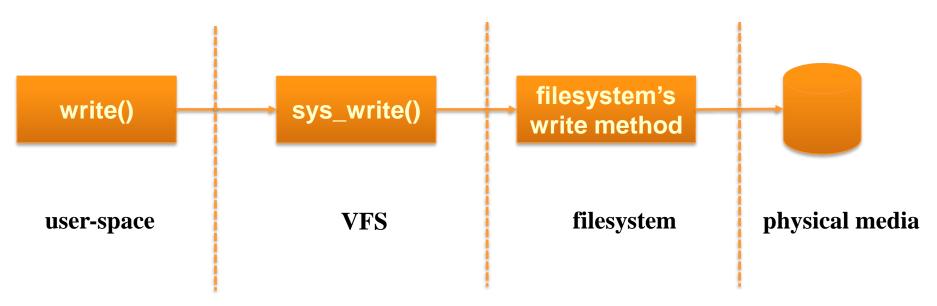
- **NO Lecture on Tuesday, November 6**

# *Lecture Topics*

- System Calls

# *System Calls (Definition & An Example)*

- Unix/Linux uses systems calls to implement most interfaces between User Mode processes and hardware device

➢ **write (fd, &buf, len)**

| write() | → | sys_write() | → | filesystem's write method | → | (disk) |
|---------|---|-------------|---|--------------------------|---|--------|

| **user-space** | **VFS** | **filesystem** | **physical media** |
|---|---|---|---|

# *Why System Call ?*

- Making programming easier (hiding low-level hardware)

- Increasing system security (checking the correctness)

- Making programs more portable (the same interfaces for different kernel)

# *System Call Categories*

– Process management

– Memory management

– File management

– Device management

– Communication

– ……

# *User Mode vs. Kernel Mode*

An interrupt or exception (INT)

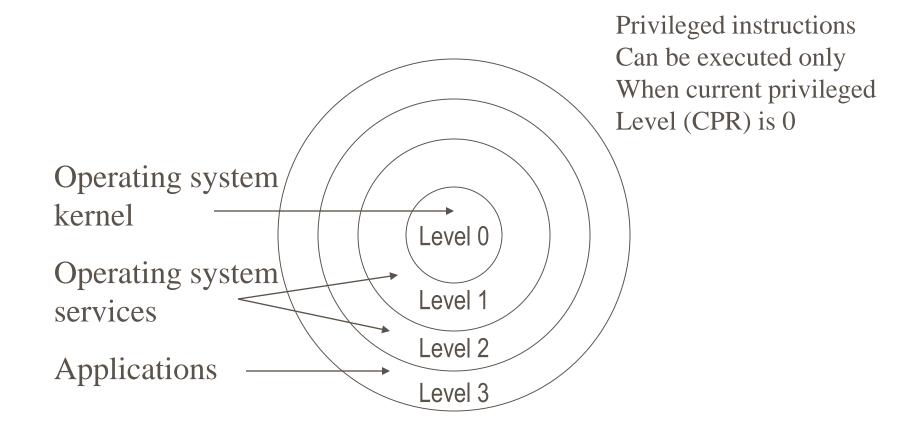| User mode | Kernel (privileged) mode |
|---|---|
| ➤Access user-mode memory | ➤Access kernel/user-mode memory |

A special instruction (IRET)
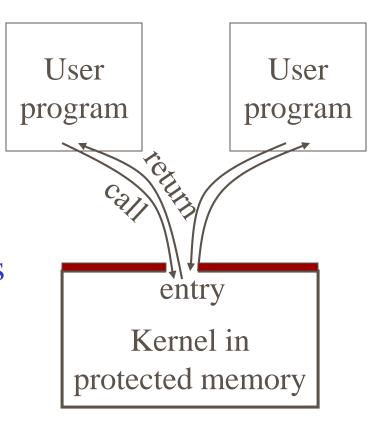
# *Why Privilege Mode ?*

- ## Special Instructions
  - Mapping, TLB, etc
  - Device registers
  - I/O channels, etc.

- ## Processor Features
  - SSE3 (Streaming Single-Instruction-Multiple-Data Extension 3)
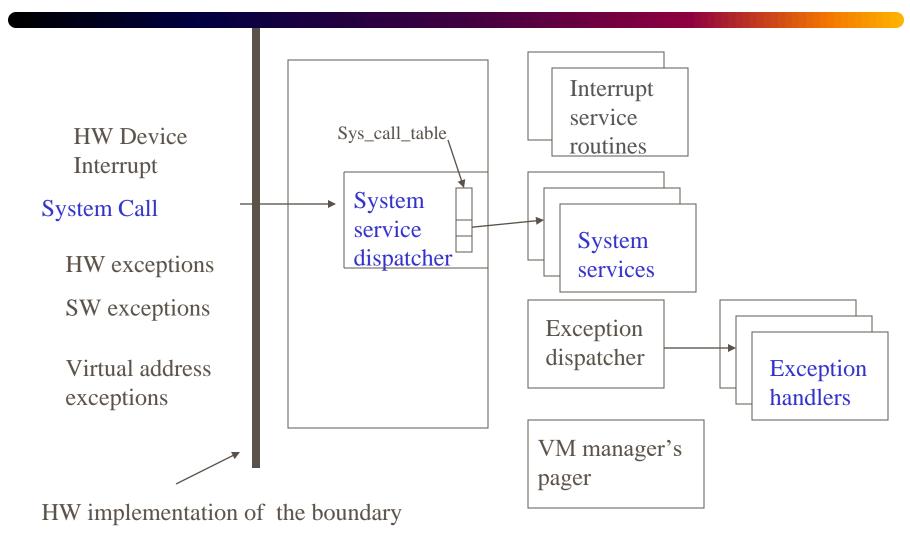
- ## Device access

# *X86 Protection Ring*

Privileged instructions
Can be executed only
When current privileged
Level (CPR) is 0

Operating system
kernel

Operating system
services

Applications

Level 0

Level 1

Level 2

Level 3

# *System Call Mechanism*

- ➤ User code can be arbitrary
- ➤ User code cannot modify kernel memory
- ➤ Makes a system call with parameters
- ➤ The call mechanism switches code to kernel mode
- ➤ Execute system call
- ➤ Return with results

User program

User program

return

call

entry

Kernel in protected memory

# System Call Implementation

HW Device
Interrupt

System Call

HW exceptions

SW exceptions

Virtual address
exceptions

HW implementation of the boundary

Sys_call_table

System
service
dispatcher

Interrupt
service
routines

System
services

Exception
dispatcher

Exception
handlers

VM manager's
pager

# *System Call Entry Point*

■ Assume passing parameters in registers

  EntryPoint:

      switch to kernel stack

      save context

      check register

      call the real code pointed by register

      restore context

      switch to user stack

      iret (change to user mode and return)

| User stack | User memory |
|---|---|
| | Registers |

| Kernel stack | Registers |
|---|---|
| | Kernel memory |

# *System Call in Linux*

- In Linux, all system calls use the following conventions
  - INT 0x80 to invoke (IDT vector 0x80)
  - EAX = system call # (asm/unistd.h)
  - EAX = return value (negative for errors) (asm-generic/errno.h)

# *System Call*

- Vector in IDT is `system_call` (in arch/i386/kernel/entry.S)

  – saves registers to stack

  – check for a valid system call #

  – call specific routine using a jump table:
    `sys_call_table` (syscall_table.S, included from entry.S)

  – stack (in kernel level) now appears as shown…
    (note the argument order)

**ESP**

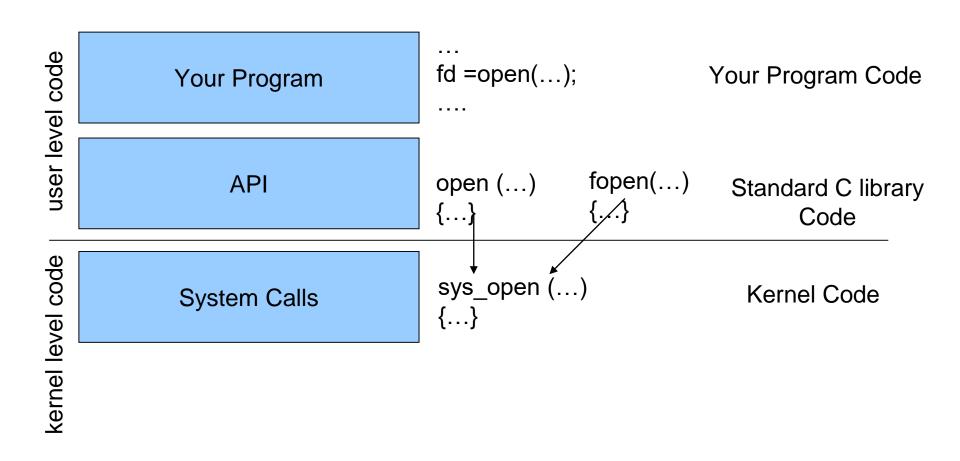| ret. addr |
| --- |
| orig. EBX |
| orig. ECX |
| orig. EDX |
| orig. ESI |
| orig. EDI |
| orig. EBP |

- Each system call is assigned a unique syscall number

```
574  .data
575  ENTRY(sys_call_table)
576        .long sys_restart_syscall      /* 0 - old "setup()" system call, used for restarting */
577        .long sys_exit
578        .long sys_fork
579        .long sys_read
580        .long sys_write
581        .long sys_open            /* 5 */
582        .long sys_close
583        .long sys_waitpid
584        .long sys_creat
585        .long sys_link
586        .long sys_unlink          /* 10 */
587        .long sys_execve
588        .long sys_chdir
589        .long sys_time
590        .long sys_mknod
591        .long sys_chmod           /* 15 */
592        .long sys_lchown16
593        .long sys_ni_syscall      /* old break syscall holder */
594        .long sys_stat
595        .long sys_lseek
596        .long sys_getpid          /* 20 */
597        .long sys_mount
598        .long sys_oldumount
599        .long sys_setuid16
600        .long sys_getuid16
```

# The Internal of A Simple System Call

Modifier for syscalls

Return type

Naming convention: sys_xxx()

```
asmlinkage    long    sys_getpid (void)
{
        return current->pid;
}
```

# *An Example: open()*

user level code

kernel level code

| Your Program | … <br> fd =open(…); <br> …. | Your Program Code |

| API | open (…) {…}     fopen(…) {…} | Standard C library Code |

| System Calls | sys_open (…) {…} | Kernel Code |

# *An Example: open()*

# *From Lib Code to System Call*

- C library code

  - arranges the arguments

  - before performing the system call

  - for example, open and sys_open

```
int open (const char* name, int flags, int mode);


asmlinkage long sys_open (const char* name,
                          int flags, int mode);
```

# *Open System Call Wrapper*

```
open: pushl     %ebx                    # save EBX to stack   collee saved register
      movl      0x10(%esp),%edx         # EDX   mode
      movl      0x0C(%esp),%ecx         # ECX   flags
      movl      0x08(%esp),%ebx         # EBX   name
      movl      $0x05,%eax              # open is system call #5
      int       $0x80                   # do the system call
      cmpl      $0xFFFFF001,%eax        # -1 to -4095 are errors
      jb        done                    # others are valid descriptors
      xorl      %edx,%edx               # negate error number
      subl      %eax,%edx
      pushl     %edx                    # save EDX (caller-saved)
      call      __errno_location        # get pointer to errno
      popl      %ecx                    # pop error number into ECX
      movl      %ecx,(%eax)             # save error number in errno
      orl       $0xFFFFFFFF,%eax        # return -1
done: popl      %ebx                    # restore EBX from stack
      ret
```

# *Open System Call Wrapper*

- Call to `open` translates to a call to `sys_open` inside the kernel
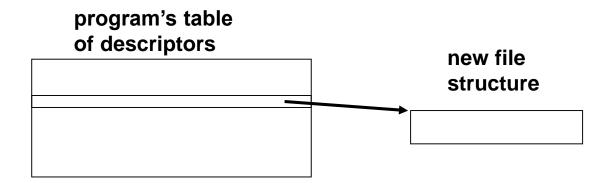

- Before going on, let's see error handling path

- Recall that, on error, C library
  - returns -1 (for many calls, including `open`)
  - stores error code in `errno` variable


- But library code is relocatable, so `errno` address is not fixed

# *Open System Call Wrapper*

- **In dynamically-linked code**
  - code addresses generally linked through jump tables
  - data addresses often use the trick
    - make a "fake" call
    - call pushes return address onto stack
    - return address is located at fixed offset from static variable
    - load from stack and add offset to obtain pointer to variable

- **`sys_open`** (in fs/open.c) invokes **`do_sys_open`**

- **`do_sys_open`** does the following

  – get a free file descriptor

  – open the file (using **`do_filp_open`**)

  – attach the file to the descriptor

**program's table
of descriptors**

**new file
structure**

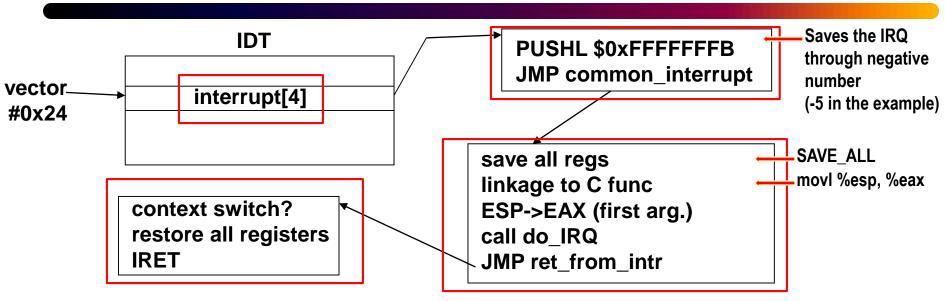# *Open System Call Wrapper*

- **`filp_open`**

  - check access rights

  - find VFS mount

  - get a file structure from a free list

  - open directory entry (using **`__dentry_open`**, below)

- **`__dentry_open`**

  - fill in file structure

  - if **`fops`** structure is not NULL && **`fops->open`** entry is not NULL, call **`fops->open`**

# *Interrupt Invocation*

**IDT**

vector #0x24 → interrupt[4]

PUSHL $0xFFFFFFFB
JMP common_interrupt

Saves the IRQ through negative number
(-5 in the example)

save all regs
linkage to C func
ESP->EAX (first arg.)
call do_IRQ
JMP ret_from_intr

SAVE_ALL
movl %esp, %eax

context switch?
restore all registers
IRET

- Why funnel all IRQs into one piece of code(instead of producing one function per IRQ) ?
    - kernel code/size versus speed tradeoff
        - perhaps no big deal for 16
        - keep in mind that you're saving a tiny number of instructions
        - APIC used by most SMPs has 256 vectors

ECE391

**After call to do_IRQ** ➡️

# *Stack on Entry to do_IRQ*

ESP → linkage return addr.

saved by SAVE_ALL

orig_eax

return address for IRET

Automatically saved by the control unit