# Chapter 2

## Operating-System Structures

Yunmin Go

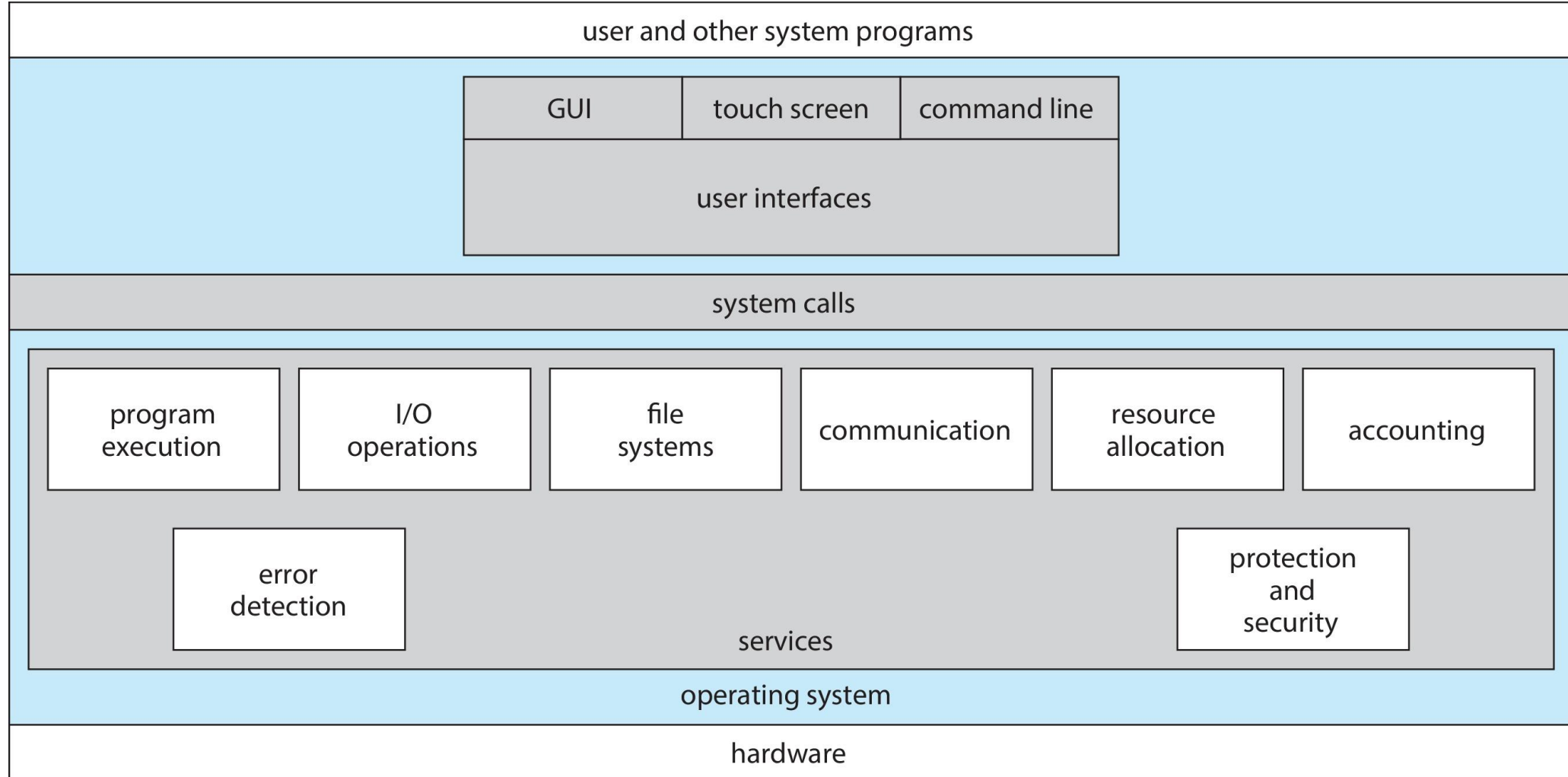School of CSEE

HANDONG GLOBAL UNIVERSITY

# Agenda

- **Operating System Services**
- System Calls
- Operating System Operations
- Operating System Structure

# Operating System Services

- Operating systems provide an environment for execution of programs and services to programs and users

- Services for user
  - User interface
  - Program execution
  - I/O operations
  - File-system manipulation
  - Communications
  - Error detection

- Functions for efficient operation of system itself
  - Resource allocation
  - Logging
  - Protection and security
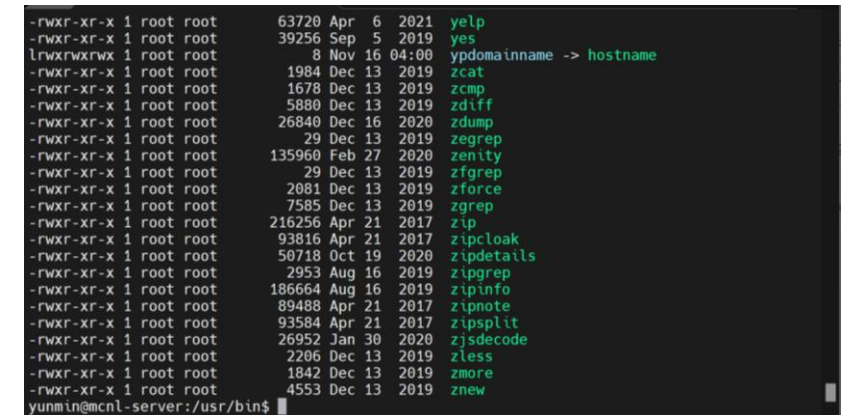
# Operating System Services

# User Interfaces: CLI

- Command line interface or command interpreter allows direct command entry
  - Sometimes implemented in kernel, sometimes by systems program
  - Sometimes multiple flavors implemented – shells
  - Primarily fetches a command from user and executes it
  - Sometimes commands built-in, sometimes just names of programs



&lt;Windows Command Prompt&gt;



&lt;Linux bash&gt;

# User Interfaces: GUI

- User-friendly desktop metaphor interface

  - Usually mouse, keyboard, and monitor

  - Icons represent files, programs, actions, etc

  - Various mouse buttons over objects in the interface cause various actions (provide information, options, execute function, open directory (known as a folder)

  - Invented at Xerox PARC

# User Interfaces: Touchscreen Interfaces

- Touchscreen devices require new interfaces
  - Mouse not possible or not desired
  - Actions and selection based on gestures
  - Virtual keyboard for text entry
- Voice commands (+AI)



https://support.apple.com/ko-kr/siri



https://medium.com/dsgnrs/iphone-x-ui-guidelines-screen-details-and-layout-9ecb795c205a

HANDONG GLOBAL UNIVERSITY

# Agenda

- Operating System Services
- **System Calls**
- Operating System Operations
- Operating System Structure

# System Calls

- System calls provide an interface to the services made available by an operating system
  - "Function calls to OS kernel available through interrupt"
  - Generally, provided as interrupt handlers written in C/C++ or assembly.
  - A mechanism to transfer control safely from lesser privileged modes to higher privileged modes.
    - Ex) POSIX system calls: open, close, read, write, fork, kill, wait, …
  - Mostly accessed by programs via a high-level Application Programming Interface (API) rather than direct system call use

# Example of System Calls

- open()

```
#include <fcntl.h>
int open(const char *pathname, int flags);
int open(const char *pathname, int flags, mode_t mode);
```

Open and possibly create a file

- Open the file specified by pathname

- If the specified file does not exist, it may optionally be created

- Return value
  - Success: file descriptor (nonnegative integer)
  - Error: -1
    - *errno* is set to indicate error

- flags
  - it must include one of access modes: O_RDONLY, O_WRONLY, and O_RDWR
  - zero or more file creation flags and file status flags can be bitwise-or'd in flags.
    - O_CREAT, O_TRUNC, O_DIRECTORY,..

- mode
  - specifies the file mode bits (permission) to be applied when a new file is created

# Example of System Calls

- ## read()

```
#include <unistd.h>
ssize_t read(int fd, void *buf, size_t count);
```

## Read from a file descriptor

- Reads up to *count* bytes from file descriptor *fd* into the buffer starting at *buf*

- Return value
  - Success: the number of bytes read
    - file position is advanced by this number
  - Fail: -1
    - *errno* is set to indicate error

HANDONG GLOBAL UNIVERSITY

# Example of System Calls

- write()

```
#include <unistd.h>
ssize_t write(int fd, const void *buf, size_t count);
```

## Write to a file descriptor

- Writes up to *count* bytes from the buffer starting at *buf* to the file referred to by the file descriptor *fd*

- Return value
  - Success: the number of bytes written is returned
    - file position is advanced by this number
  - Fail: -1
    - *errno* is set to indicate error

HANDONG GLOBAL UNIVERSITY

# Example of System Calls

- System call sequence to print the contents of file

```c
int main(int argc, char *argv[]) {
    int fd;
    int start_offset;
    int read_bytes;
    char filename[FILELEN];
    char buf[BUFSIZE];

    if (argc < 3) {
        fprintf(stderr ,"Usage: %s [file name]
                [start offset] \n", argv[0]);
        exit(1);
    }

    strcpy(filename, argv[1]);      // set file name
    start_offset = atoi(argv[2]);   // set start offset
```

```c
    // open source file
    fd = open(filename, O_RDONLY);
    if (fd == -1) {
        perror(filename);
        exit(1);
    }

    // set file position
    lseek(fd, start_offset, SEEK_SET);

    // read data from file
    while ((read_bytes = read(fd, buf, BUFSIZE)) > 0) {
        // print file data
        write(1, buf, read_bytes);
    }
    printf("\n\n");

    // close file
    close(fd);

    return 0;
}
```

# Example of System Calls

- System call sequence to print the contents of file
  - Results

```
yunmin@mcn1-server:~/workspace/os/ch2$ gcc printfile.c -o printfile
yunmin@mcn1-server:~/workspace/os/ch2$ ./printfile
Usage: ./printfile [file name] [start offset]
yunmin@mcn1-server:~/workspace/os/ch2$ ./printfile hello.c 0
#include <stdio.h>

int main()
{
    printf("Hello Handong!\n");
    return 0;
}

yunmin@mcn1-server:~/workspace/os/ch2$ ./printfile hello.c 10
stdio.h>

int main()
{
    printf("Hello Handong!\n");
    return 0;
}
```

# Types of System Calls

- Process control
- File management
- Device management
- Information maintenance
- Communications
- Protection

| | Windows | Unix |
|---|---|---|
| Process control | CreateProcess()<br>ExitProcess()<br>WaitForSingleObject() | fork()<br>exit()<br>wait() |
| File management | CreateFile()<br>ReadFile()<br>WriteFile()<br>CloseHandle() | open()<br>read()<br>write()<br>close() |
| Device management | SetConsoleMode()<br>ReadConsole()<br>WriteConsole() | ioctl()<br>read()<br>write() |
| Information maintenance | GetCurrentProcessID()<br>SetTimer()<br>Sleep() | getpid()<br>alarm()<br>sleep() |
| Communications | CreatePipe()<br>CreateFileMapping()<br>MapViewOfFile() | pipe()<br>shm_open()<br>mmap() |
| Protection | SetFileSecurity()<br>InitlializeSecurityDescriptor()<br>SetSecurityDescriptorGroup() | chmod()<br>umask()<br>chown() |

# System Programs

- System programs provide a convenient environment for program development and execution. They can be divided into:
  - File manipulation: create, delete, copy, rename, print, dump, list, and generally manipulate files and directories
  - Status information: date, time, system information, performance, etc.
  - Programming language support: compiler, assembler, debuggers, interpreter
  - Program loading and execution: loader, linkage editors, etc.
  - Communications: networks
  - Background services: launch at boot time, process scheduling, logging, etc
  - Application programs: run by users
- Most users' view of the operation system is defined by system programs, not the actual system calls

HANDONG GLOBAL UNIVERSITY

# Why Applications are Operating System Specific

- Apps compiled on one system usually not executable on other operating systems

- Each operating system provides its own unique system calls
  - Own file formats, etc

- Apps can be multi-operating system
  - Written in interpreted language like Python, Ruby, and interpreter available on multiple operating systems
  - App written in language that includes a VM containing the running app (like Java)
  - Use standard language (like C), compile separately on each operating system to run on each

- Application Binary Interface (ABI) is architecture equivalent of API, defines how different components of binary code can interface for a given operating system on a given architecture, CPU, etc

# Agenda

- Operating System Services
- System Calls
- **Operating System Operations**
- Operating System Structure

# Operating-System Operations

- Bootstrap program
  - simple code to initialize the system, load the kernel
- Kernel loads
- Starts system daemons (services provided outside of the kernel)
- Kernel interrupt driven (hardware and software)
  - Hardware interrupt by one of the devices
  - Software interrupt (exception or trap):
    - Software error (e.g., division by zero)
    - Request for operating system service: system call
    - Other process problems include infinite loop, processes modifying each other or the operating system

HANDONG GLOBAL UNIVERSITY

# Multiprogramming

- **Multiprogramming** needed for efficiency
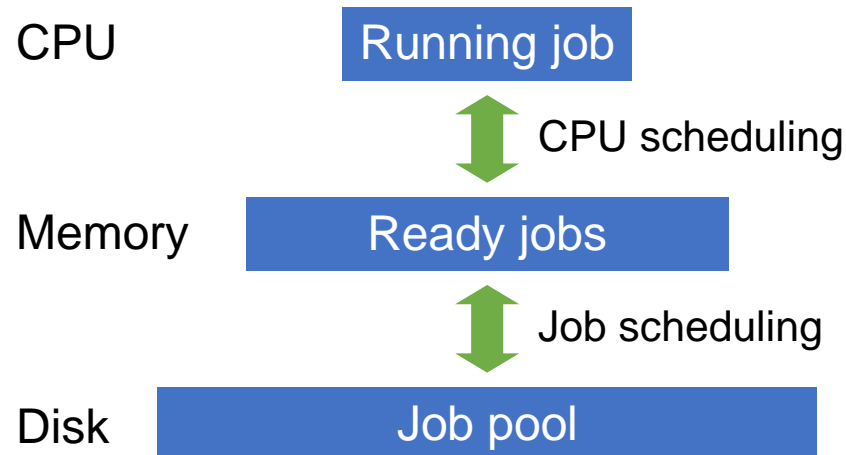  - Single user cannot keep CPU and I/O devices busy at all times
  - Multiprogramming organizes jobs (code and data) so CPU always has one to execute
  - A subset of total jobs in system is kept in memory
  - One job selected and run via job scheduling
  - When it has to wait (for I/O for example), OS switches to another job



<Memory layout for a multiprogramming system>

# Multitasking

- Multitasking is logical extension of multiprogramming
  - CPU switches jobs so frequently that users can interact with each job while it is running, creating interactive computing
  - Response time should be < 1 second
  - Each user has at least one program executing(process) in memory
  - CPU scheduler selects a job that is ready to run

CPU    | Running job |

↕ CPU scheduling

Memory | Ready jobs |

↕ Job scheduling

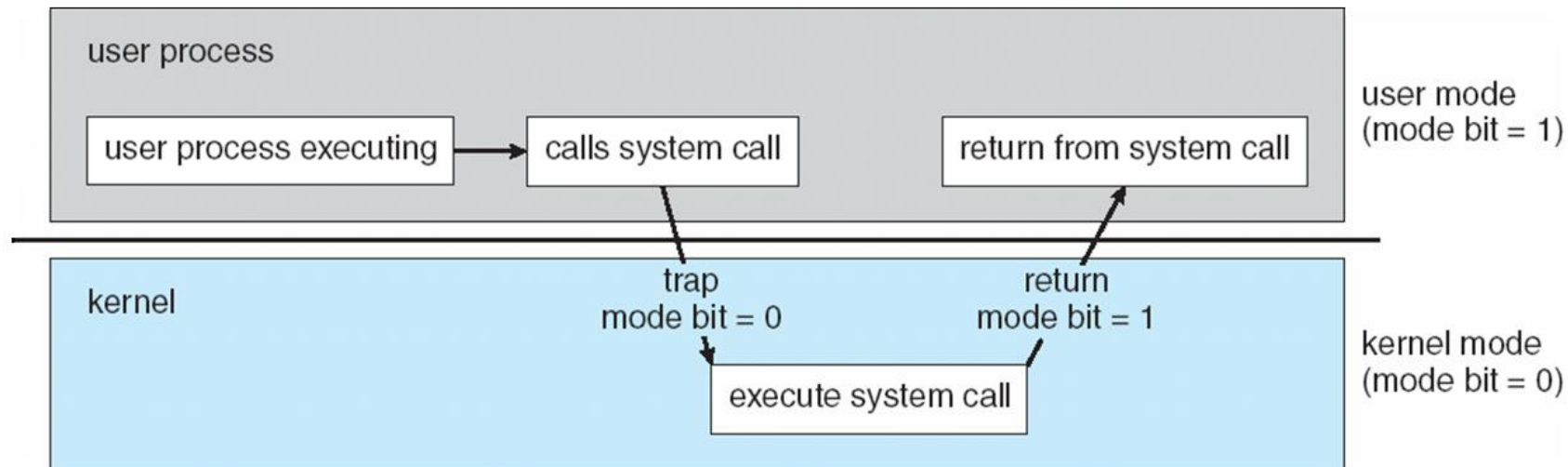Disk   | Job pool |

HANDONG GLOBAL UNIVERSITY

# Requirement for Modern OS

- Modern OS's are interrupt-driven programs
    - Events are signaled by interrupts, which are handled by interrupt handlers.
    - H/W and S/W resources are shared.
    - ➡ **Problem: An error from a process can corrupt whole system**

- Essential requirement for multi-user OS: Error of a program should not affect to other program
    - Dangerous instructions
    - Preventing long execution of user process

# Dual-Mode Operation

- Dual-mode operation allows OS to protect itself and other system components (requires H/W support)
  - User mode
    - User defined code (application)
    - Privileged instructions, which can cause harm to other system, are prohibited
  - Kernel mode (supervisor mode, system mode, privileged mode)
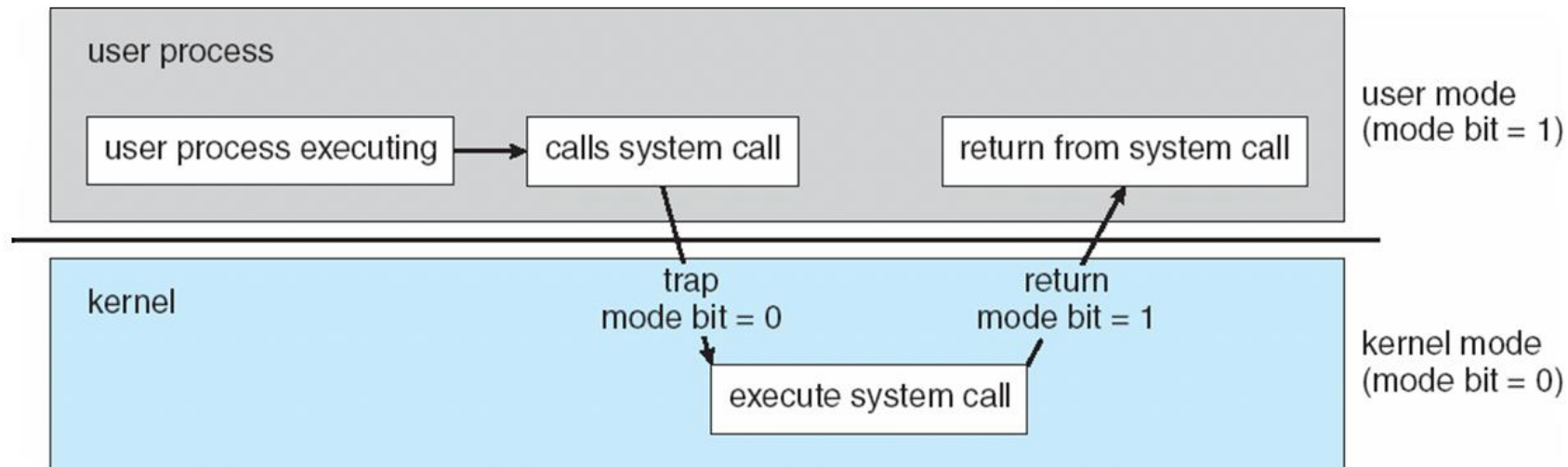    - OS code, System call → Privileged instructions are permitted
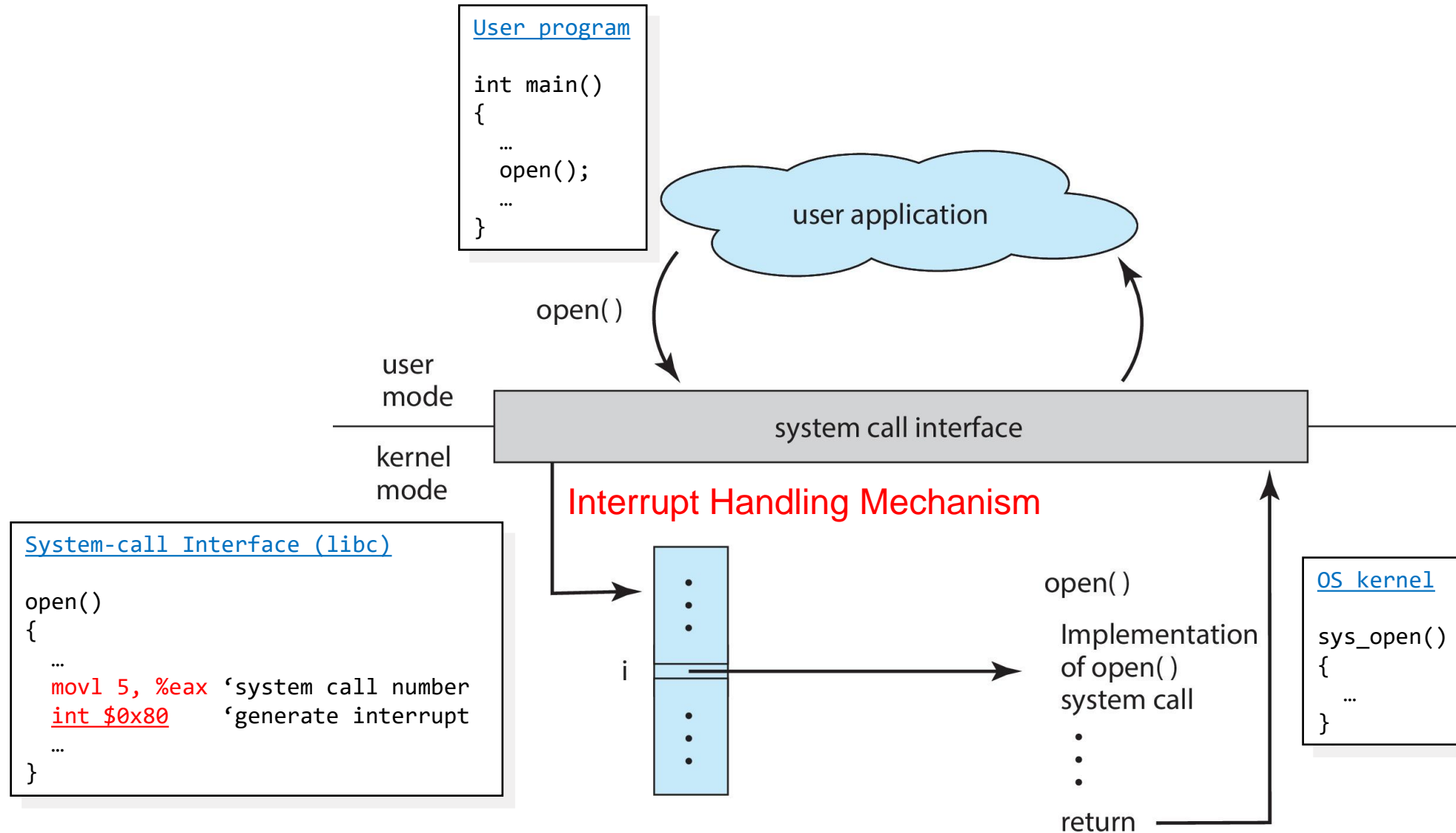
# Dual-Mode Operation

- Dual-mode operation allows OS to protect itself and other system components (requires H/W support)
  - Mode bit provided by hardware
    - Provides ability to distinguish when system is running user code or kernel code.
    - When a user is running → mode bit is "user" (1)
    - When kernel code is executing → mode bit is "kernel" (0)

# System Call – OS Relationship



User program

```
int main()
{
    …
    open();
    …
}
```

user application

open()

user
mode

kernel
mode

system call interface

Interrupt Handling Mechanism

System-call Interface (libc)

```
open()
{
    …
    movl 5, %eax 'system call number
    int $0x80    'generate interrupt
    …
}
```

i

open()

Implementation
of open()
system call

return

OS kernel
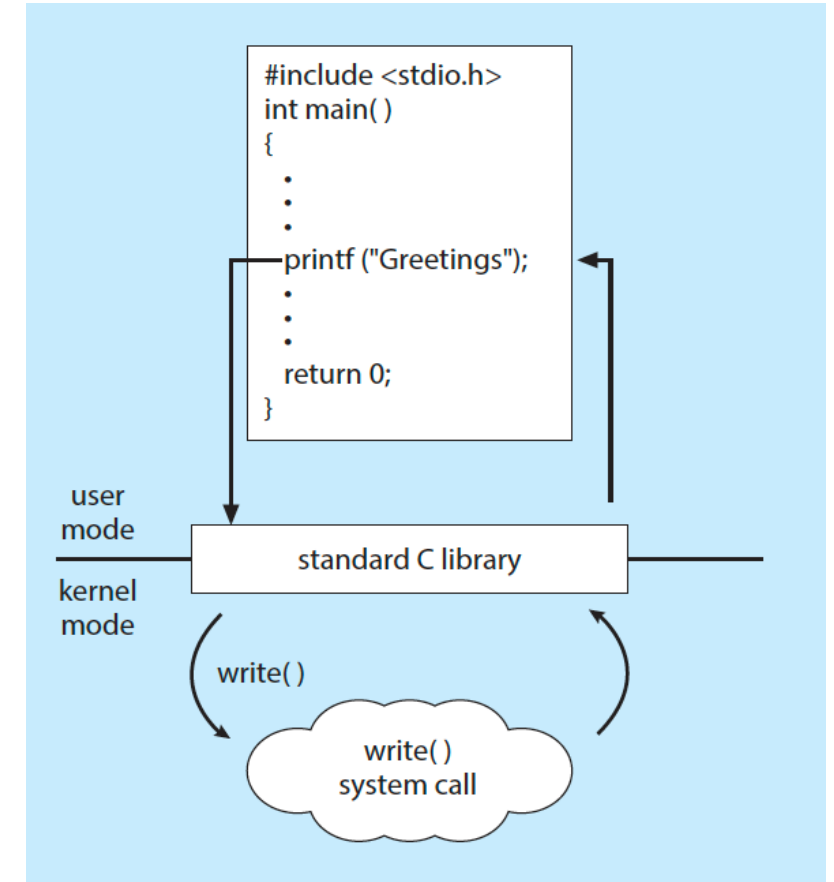
```
sys_open()
{
    …
}
```

HANDONG GLOBAL
UNIVERSITY

# System Call Implementation

- Typically, a number associated with each system call
  - System-call interface maintains a table indexed according to these numbers
- The system call interface invokes the intended system call in OS kernel and returns status of the system call and any return values
  - Passing information to the kernel
  - Switch to kernel mode
  - Any data processing and preparation for execution in kernel mode
- The caller need know nothing about how the system call is implemented
  - Just needs to obey API and understand what OS will do as a result call
  - Most details of OS interface hidden from programmer by API

HANDONG GLOBAL UNIVERSITY

# System Call Implementation

- System call vs. I/O functions in programming language

  - write() vs. printf()

    - printf() is implemented using write()

      - write(): provided by OS
      - printf(): standard function defined in C language

  - read() vs. fread()

    - fread() is implemented using read()

# System Call Implementation

- System call vs. I/O functions in programming language
  - `strace`: traces system calls invoked by a process

```
yunmin@mcnl-alpha:~/lecture/os/ch02$ strace ./hello
execve("./hello", ["./hello"], 0x7ffcfbe05690 /* 44 vars */) = 0
brk(NULL)                               = 0x5a87669c0000
arch_prctl(0x3001 /* ARCH_??? */, 0x7fffc48148d0) = -1 EINVAL (Invalid argument)
mmap(NULL, 8192, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_ANONYMOUS, -1, 0) = 0x7786021fa000
access("/etc/ld.so.preload", R_OK)      = -1 ENOENT (No such file or directory)
openat(AT_FDCWD, "/etc/ld.so.cache", O_RDONLY|O_CLOEXEC) = 3
newfstatat(3, "", {st_mode=S_IFREG|0644, st_size=60255, ...}, AT_EMPTY_PATH) = 0
mmap(NULL, 60255, PROT_READ, MAP_PRIVATE, 3, 0) = 0x7786021eb000

.......

getrandom("\x63\x6a\x48\x81\x43\x3c\xcf\x2f", 8, GRND_NONBLOCK) = 8
brk(NULL)                               = 0x5a87669c0000
brk(0x5a87669e1000)                     = 0x5a87669e1000
write(1, "Hello Handong!\n", 15Hello Handong!
)       = 15
exit_group(0)                           = ?
+++ exited with 0 +++
```

# System Call Implementation
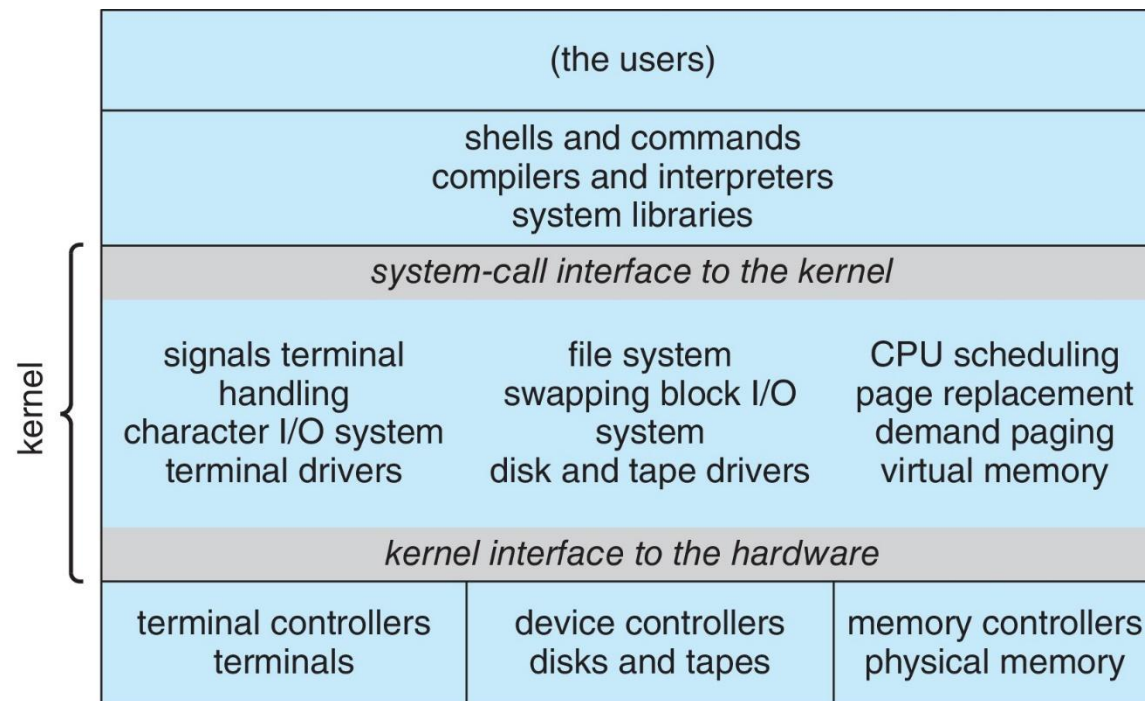
- Let's try running gdb!

# Agenda

- Operating System Services
- System Calls
- Operating System Operations
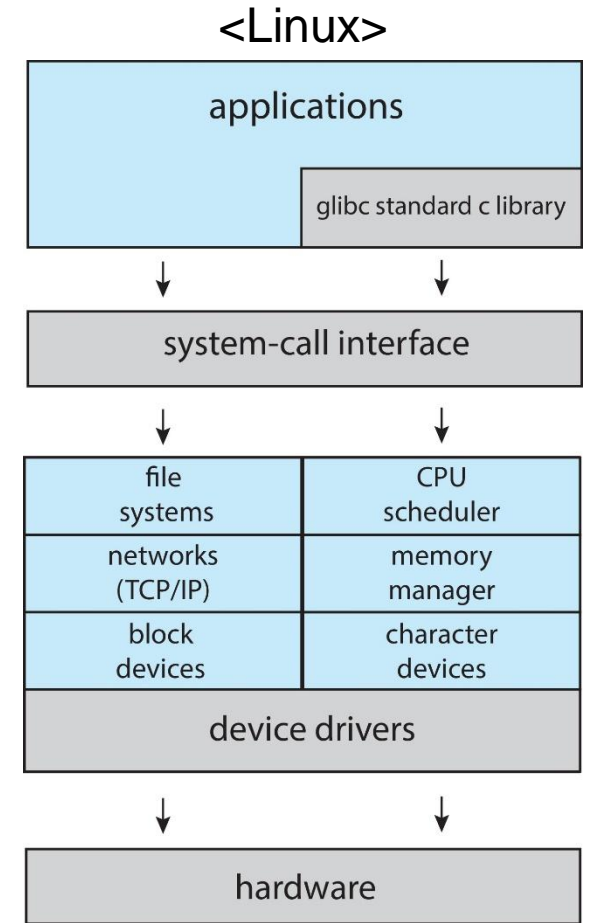- **Operating System Structures**

# Monolithic Structure

- Monolithic structure place all of the functionality of the kernel into a single, static library file that runs in a single address space.

<Traditional UNIX System>

| (the users) | | |
|---|---|---|
| shells and commands<br>compilers and interpreters<br>system libraries | | |
| *system-call interface to the kernel* | | |
| signals terminal<br>handling<br>character I/O system<br>terminal drivers | file system<br>swapping block I/O<br>system<br>disk and tape drivers | CPU scheduling<br>page replacement<br>demand paging<br>virtual memory |
| *kernel interface to the hardware* | | |
| terminal controllers<br>terminals | device controllers<br>disks and tapes | memory controllers<br>physical memory |

kernel

HANDONG GLOBAL UNIVERSITY

# Modular Structure

■ Many modern operating systems implement loadable kernel modules (LKM)

- Uses object-oriented approach
- Each core component is separate
- Each talks to the others over known interfaces
- Each is loadable as needed within the kernel
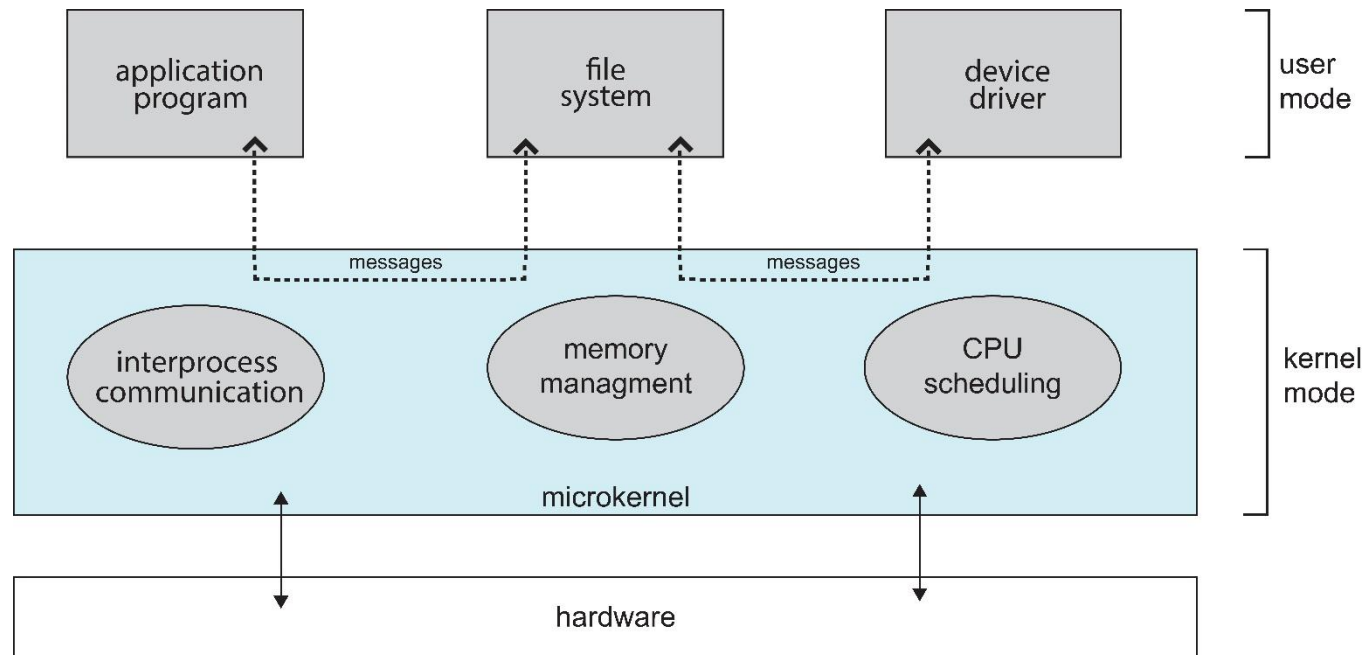- ex) Linux, Solaris, etc



\<Solaris\>



\<Linux\>

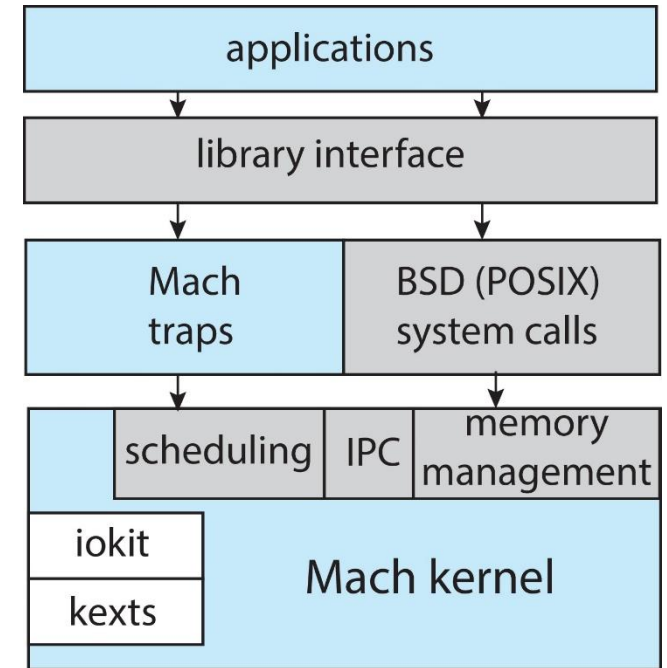Monolithic + modular design (Hybrid)

# Microkernel

- ## Smaller kernel
  - All unessential components are not implemented in kernel but as system/user-level programs.
    - Generally, process/memory management, communication facility are in the kernel.
  - System calls are provided through message passing.

# Hybrid Structure

- Most modern operating systems are actually not one pure model

  - Hybrid combines multiple approaches to address performance, security, usability needs

  - Linux and Solaris kernels in kernel address space, so monolithic, plus modular for dynamic loading of functionality

  - Windows mostly monolithic, plus microkernel for different subsystem personalities

  - macOS uses Mach microkernel with modules and integrates BSD subsystem.



<The structure of Darwin>

# Summary of Operating System Structures

| Structures | Description | Advantages | Disadvantages | Example OS |
|---|---|---|---|---|
| Monolithic | All OS services run in kernel space as a single program. | Fast execution, direct system call handling, efficient performance. | Hard to debug and maintain, a crash affects the whole system. | Linux, Unix, MS-DOS |
| Modular | A monolithic kernel with dynamically loadable modules. | Modular and flexible, allows adding/removing drivers without rebooting. | Still runs most services in kernel space, potential security issues. | Linux (with Loadable Kernel Modules), FreeBSD |
| Microkernel | Only essential services run in kernel space, while others run in user space. | Better security and stability, crashes in user space don't affect kernel. | Slower due to message passing between kernel and user space services. | Minix, QNX, seL4 |
| Hybrid | A combination of monolithic and microkernel features for performance and flexibility. | Balances performance and modularity, reduces microkernel overhead. | Increased complexity, more difficult to optimize. | macOS (XNU), Windows NT, Google Fuchsia |