

CS 318 Principles of Operating Systems

Fall 2017

Lecture 2: Architecture Support for OS

Ryan Huang

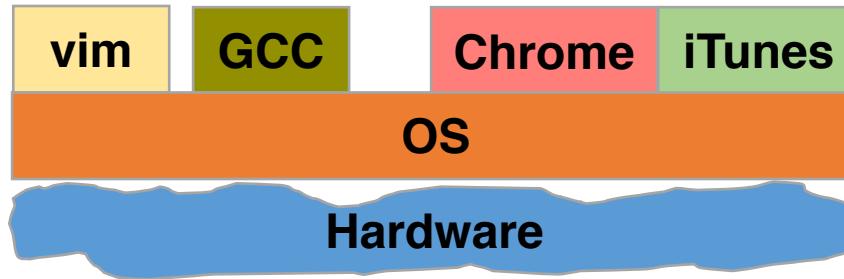


JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

Administrivia

- **Lab 0**
 - Due 09/14, done **individually**
- **Project groups**
 - Fill out Google form (link in Piazza post)
 - Just need one per group
 - Fill out even if you are working alone

Why Start With Hardware?



- **OS functionality depends upon the architectural features of H/W**
 - Key goals of an OS are to enforce [protection](#) and [resource sharing](#)
 - If done well, applications can be oblivious to HW details
- **Architectural support can greatly simplify or complicate OS tasks**
 - Early PC operating systems (DOS, MacOS) lacked virtual memory in part because the architecture did not support it
 - Early Sun 1 computers used two M68000 CPUs to implement virtual memory (M68000 did not have VM hardware support)

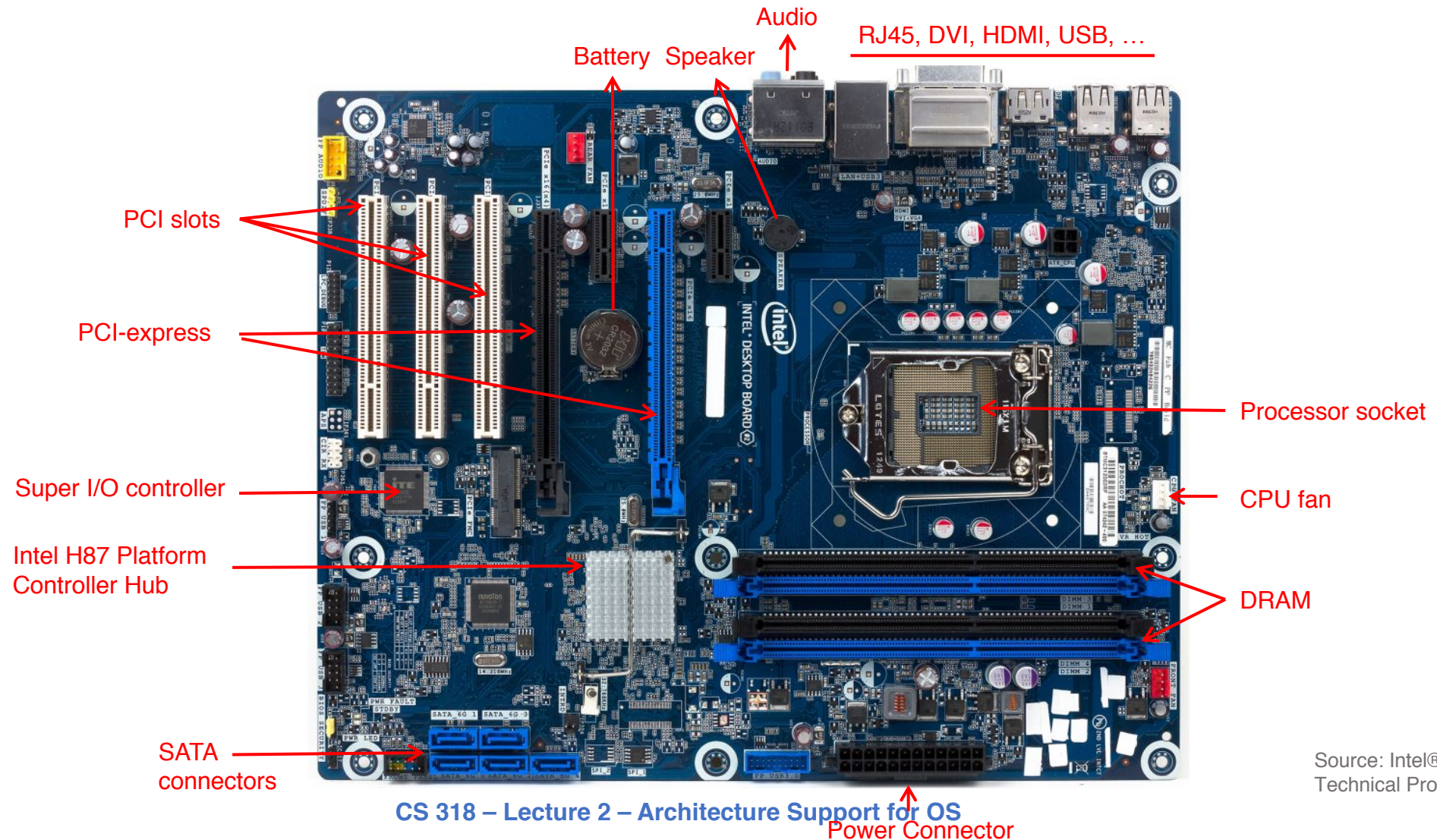
Architectural Features for OS

- **Features that directly support the OS include**
 - Bootstrapping (Lab 0)
 - Protection (kernel/user mode)
 - Protected instructions
 - Memory protection
 - System calls
 - Interrupts and exceptions
 - Timer (clock)
 - I/O control and operation
 - Synchronization

Types of Arch Support

- **Manipulating privileged machine state**
 - Protected instructions
 - Manipulate device registers, TLB entries, etc.
- **Generating and handling “events”**
 - Interrupts, exceptions, system calls, etc.
 - Respond to external events
 - CPU requires software intervention to handle fault or trap
- **Mechanisms to handle concurrency**
 - Interrupts, atomic instructions

A Motherboard (Intel DH87MC)



Source: Intel® Desktop Board DH87MC
Technical Product Specification

Typical PC System Architecture

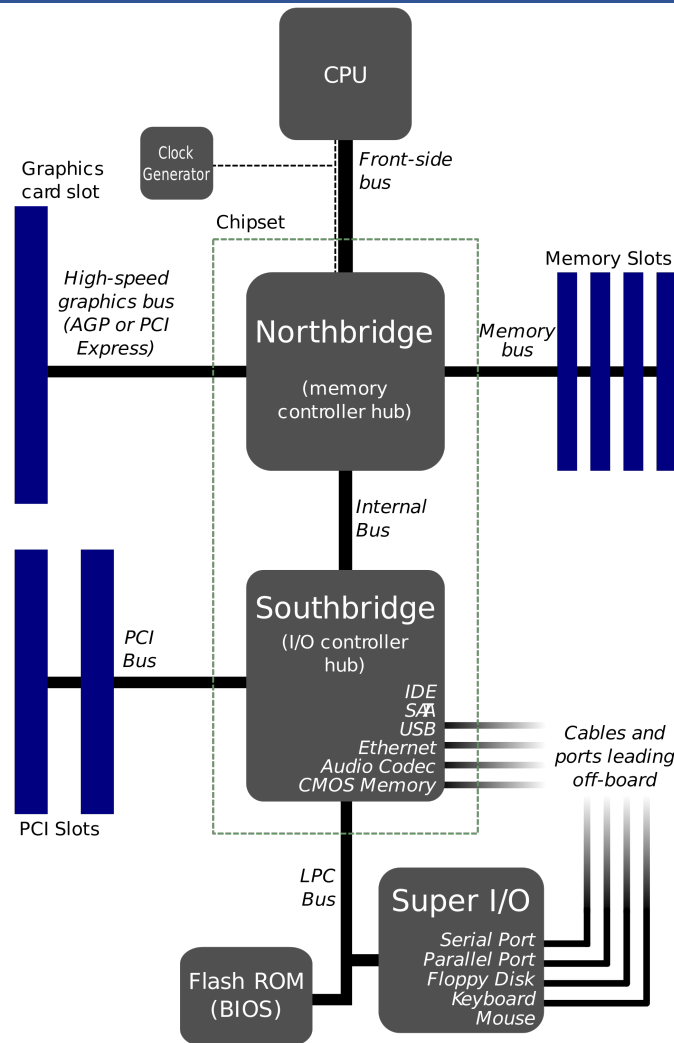


diagram source: wikipedia

Thought Experiment: A World of Anarchy

- **Any program in the system can...**
 - Directly access I/O devices
 - Write anywhere in memory
 - Read content from any memory address
 - Execute machine halt instruction
- **Do you trust such system?**
 - use Facebook
 - run Banking app
- **Challenge: protection**
 - How to execute code with restricted privilege

Thought Experiment: A Solution

- **How can we implement execution with limited privilege?**
 - Execute each program instruction through a simulator (OS)
 - If the instruction is permitted, do the instruction
 - Otherwise, stop the process
 - Basic model in Javascript and other interpreted languages
- **How do we go faster?**
 - **Most instructions are perfectly safe!**
 - Run the unprivileged code directly on the CPU
 - Do the check in h/w

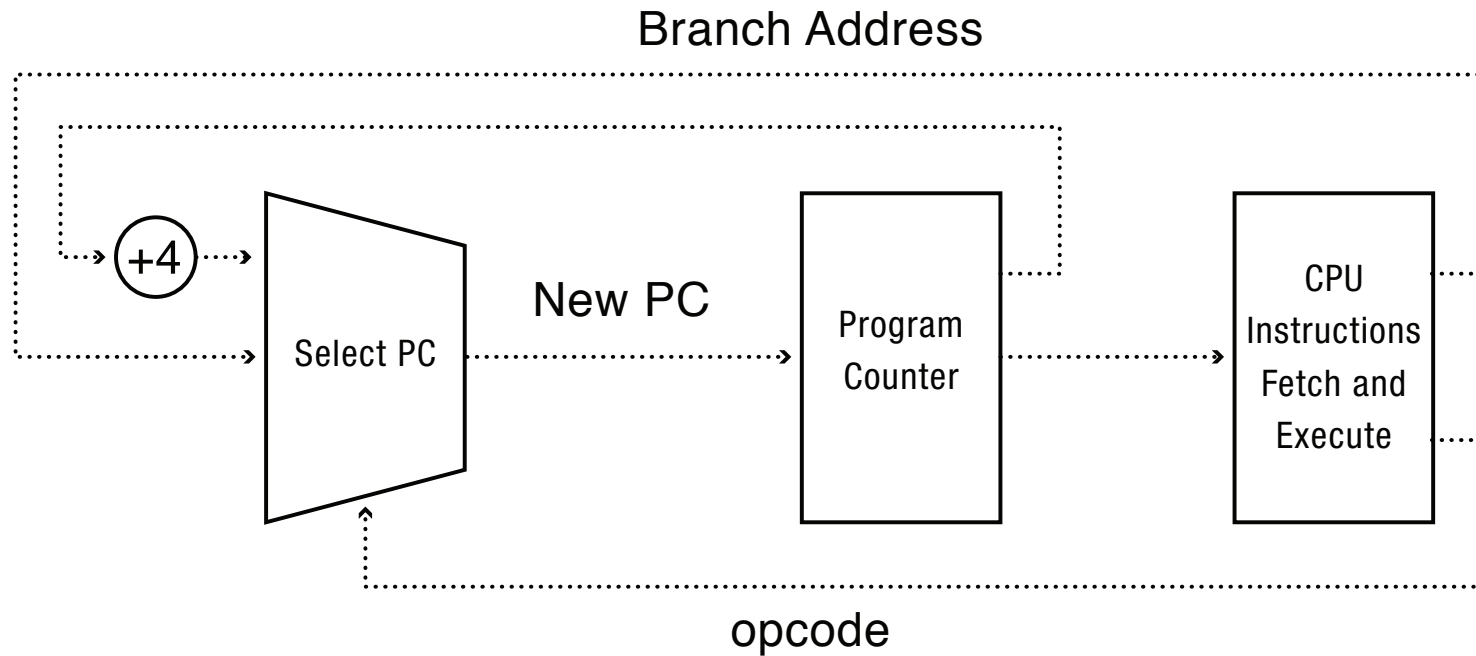
H/W Support: Dual-Mode Operation in CPU

- **User mode**
 - Limited privileges
 - Only those granted by the operating system kernel
- **Kernel mode**
 - Execution with the full privileges of the hardware
 - Read/write to any memory, access any I/O device, read/write any disk sector, send/read any packet
- **On the x86, mode stored in EFLAGS register**
- **On the MIPS, mode in the status register**

A Simple Model of a CPU

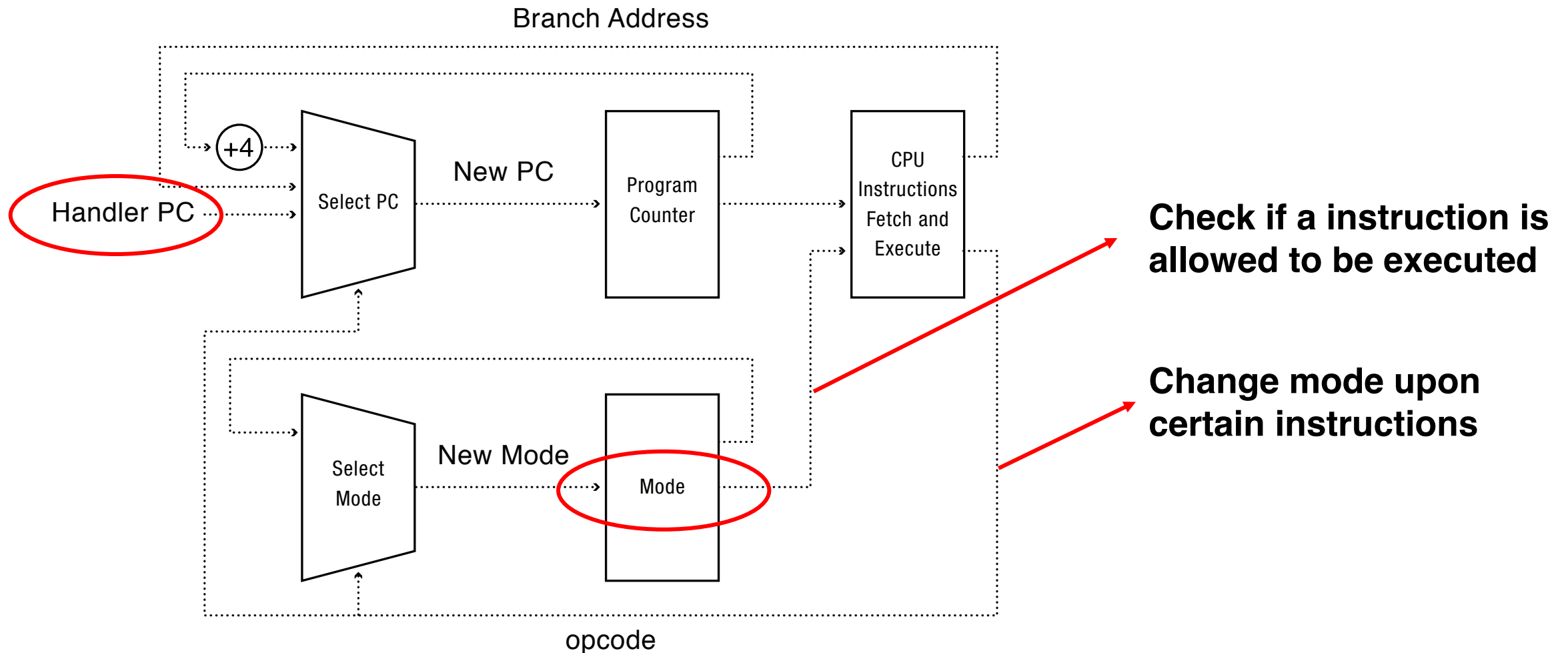
- **Basic operation**

- Fetch next instruction → Decode instruction → Execute ↪



add/subtract,
read/write memory,
call function, branch...

A CPU with Dual-Mode Operation



Protected Instructions

- **A subset of instructions restricted to use only by the OS**
 - Known as protected (privileged) instructions
- **Only the operating system can ...**
 - Directly access I/O devices (disks, printers, etc.)
 - Security, fairness (why?)
 - Manipulate memory management state
 - Page table pointers, page protection, TLB management, etc.
 - Manipulate protected control registers
 - Kernel mode, interrupt level
 - Halt instruction (why?)

INVLPG—Invalidate TLB Entries

Opcode	Instruction	Op/ En	64-Bit Mode	Compat/ Leg Mode	Description
0F 01/7	INVLPG <i>m</i>	M	Valid	Valid	Invalidate TLB entries for page containing <i>m</i> .

NOTES:

* See the IA-32 Architecture Compatibility section below.

Instruction Operand Encoding

Op/En	Operand 1	Operand 2	Operand 3	Operand 4
M	ModRM:r/m (r)	NA	NA	NA

Description

Invalidates any translation lookaside buffer (TLB) entries specified with the source operand. The source operand is a memory address. The processor determines the page that contains that address and flushes all TLB entries for that page.¹

The INVLPG instruction is a privileged instruction. When the processor is running in protected mode, the CPL must be 0 to execute this instruction.

The INVLPG instruction normally flushes TLB entries only for the specified page, however, in some cases, it may flush more entries, even the entire TLB. The instruction is guaranteed to invalidate only TLB entries associated with the current PCID. (If PCIDs are disabled — CR4.PCIDE = 0 — the current PCID is 000H.) The instruction also invalidates any global TLB entries for the specified page, regardless of PCID.

For more details on operations that flush the TLB, see “MOV—Move to/from Control Registers” and Section 4.10.4.1, “Operations that Invalidate TLBs and Paging-Structure Caches,” of the *Intel® 64 and IA-32 Architectures Software Developer’s Manual, Volume 3A*.

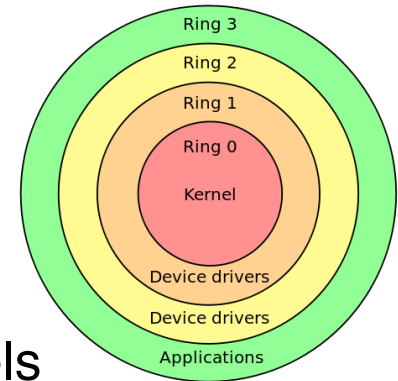
This instruction’s operation is the same in all non-64-bit modes. It also operates the same in 64-bit mode, except if the memory address is in non-canonical form. In this case, INVLPG is the same as a NOP.

IA-32 Architecture Compatibility

The INVLPG instruction is implementation dependent, and its function may be implemented differently on different

Beyond Dual-Mode Operations

- **(Modern) CPU may provide more than 2 privilege levels**
 - Called hierarchical protection domains or protection rings
 - higher-privilege rings can access lower-privilege operations and data
 - x86 supports four levels:
 - ring 0 = kernel mode, ring 3 = user mode
 - Multics provides 8 levels of privilege
 - Even seen in mobile devices
 - ARMv7 processors in modern smartphones have 8 different protection levels
- **Why?**
 - Protect the OS from itself (software engineering)
 - reserved for vendor, e.g., virtualization

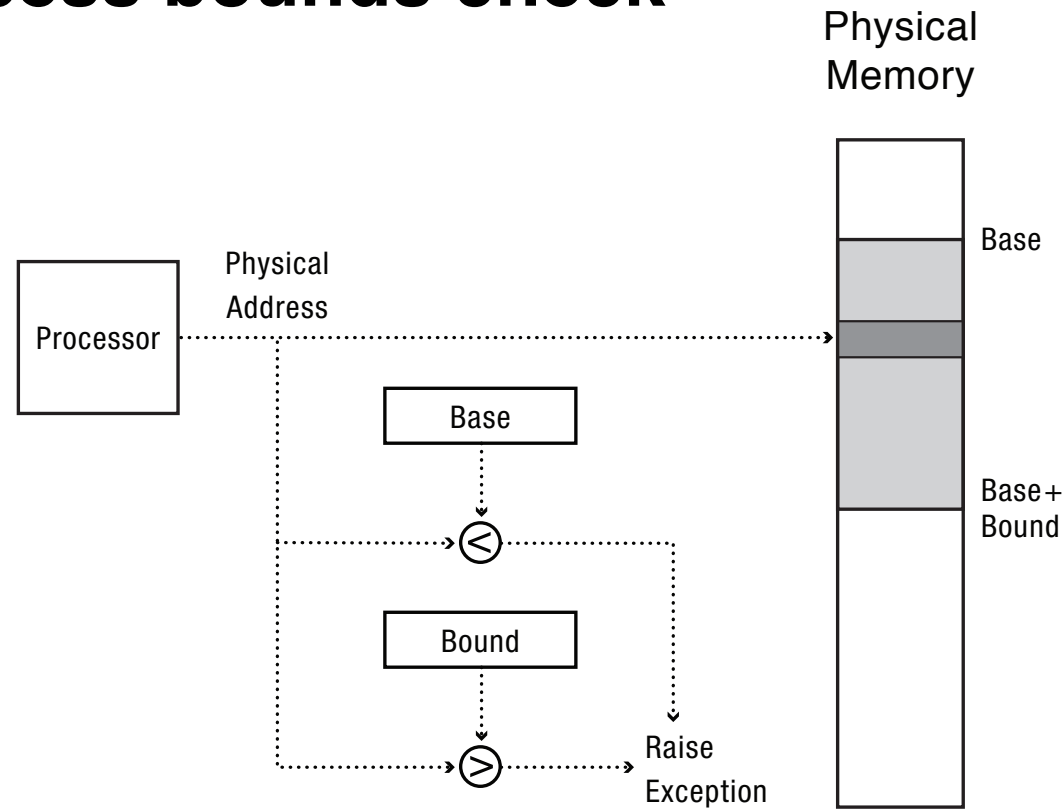


Memory Protection

- **Why?**
 - OS must be able to protect programs from each other
 - OS must protect itself from user programs
- **May or may not protect user programs from OS**
 - Raises question of whether programs should trust the OS
 - Untrusted operating systems? (Intel SGX)
- **Memory management hardware (MMU) provides the mechanisms**
 - Base and limit registers
 - Page table pointers, page protection, segmentation, TLB
 - **Manipulating the hardware uses protected (privileged) operations**

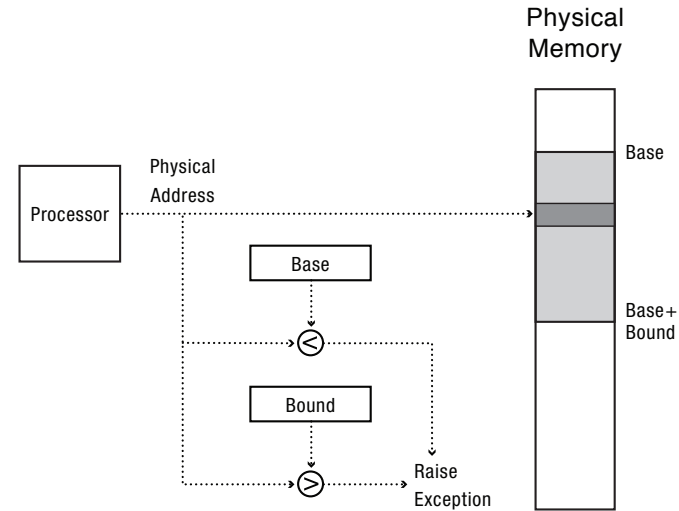
Simple Memory Protection

- **Memory access bounds check**



Simple Memory Protection

- **Memory access bounds check**
- **Problem**
 - Inflexible:
 - fix allocation, difficult to expand heap and stack
 - Memory sharing
 - Memory fragmentation
 - Physical memory
 - require changes to mem instructions each time the program is loaded

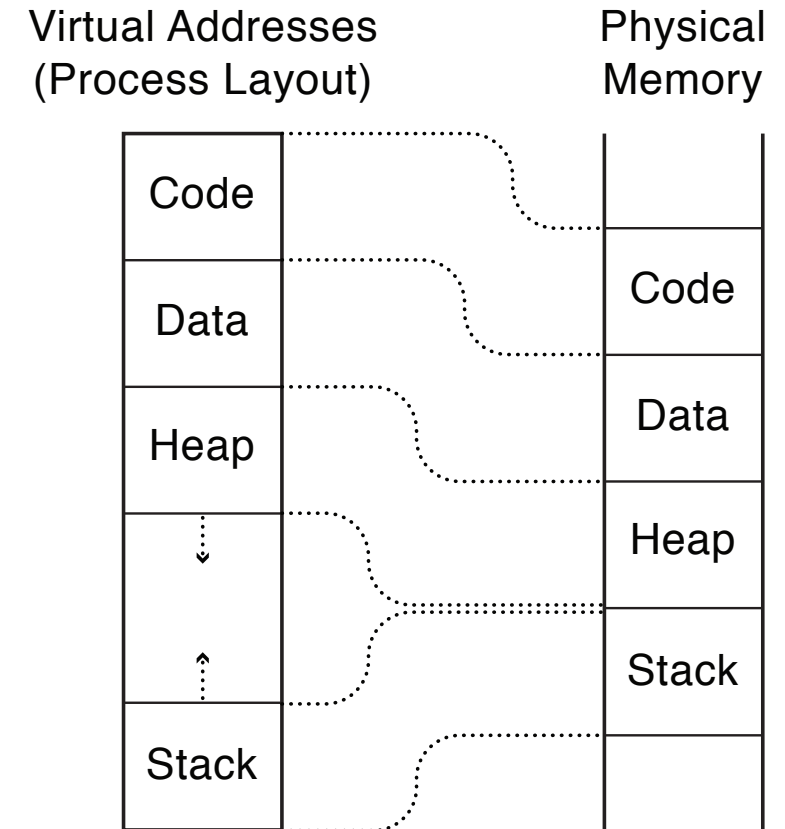


Idea: Virtual Address

- **Programs refer to memory by virtual addresses**
 - start from 0
 - illusion of “owning” the entire memory address space
- **The virtual address is translated to physical address**
 - upon each memory access
 - done in hardware using a table
 - table setup by the OS

Idea: Virtual Address

- **Programs refer to memory by virtual addresses**
 - start from 0
 - illusion of “owning” the entire memory address space
- **The virtual address is translated to physical address**
 - upon each memory access
 - done in hardware using a table
 - table setup by the OS



Types of Arch Support

- **Manipulating privileged machine state**
 - Protected instructions
 - Manipulate device registers, TLB entries, etc.
- **Generating and handling “events”**
 - Interrupts, exceptions, system calls, etc.
 - Respond to external events
 - CPU requires software intervention to handle fault or trap
- **Mechanisms to handle concurrency**
 - Interrupts, atomic instructions

OS Control Flow

- **After OS booting, all entry to kernel is a result of some event**
 - event immediately stops current execution
 - changes mode to kernel mode
 - invoke a piece of code to handle event (event handler)
- **When the processor receives an event of a given type, it**
 - transfers control to handler within the OS
 - handler saves program state (PC, regs, etc.)
 - handler executes core functionality
 - e.g., writing data to disk
 - handler restores program state, returns to program

Events

- **An event is an “unnatural” change in control flow**
 - Events immediately stop current execution
 - Changes mode, context (machine state), or both
- **The kernel defines a handler for each event type**
 - The specific types of events are defined by the architecture
 - e.g., timer event, I/O interrupt, system call trap
 - In effect, the operating system is one big event handler

Event: Interrupt vs. Exceptions

- Two kinds of events, **interrupts** and **exceptions**
- Interrupts are caused by an external event (**asynchronous**)
 - Device finishes I/O, timer expires, etc.
- Exceptions are caused by executing instructions (**synchronous**)
 - CPU requires software intervention to handle a fault or trap
 - e.g., the x86 INT instruction
 - e.g., a page fault, write to a read-only page
 - a deliberate exception is a “**trap**”, unexpected is a “**fault**”

Interrupts

- **Interrupts signal asynchronous events**
 - Indicates some device needs services
 - I/O hardware interrupts
 - Software and hardware timers
- **Why?**
 - A computer is more than CPU
 - keyboard, disk, printer, camera, etc.
 - These devices occasionally need attention
 - but can't predict when

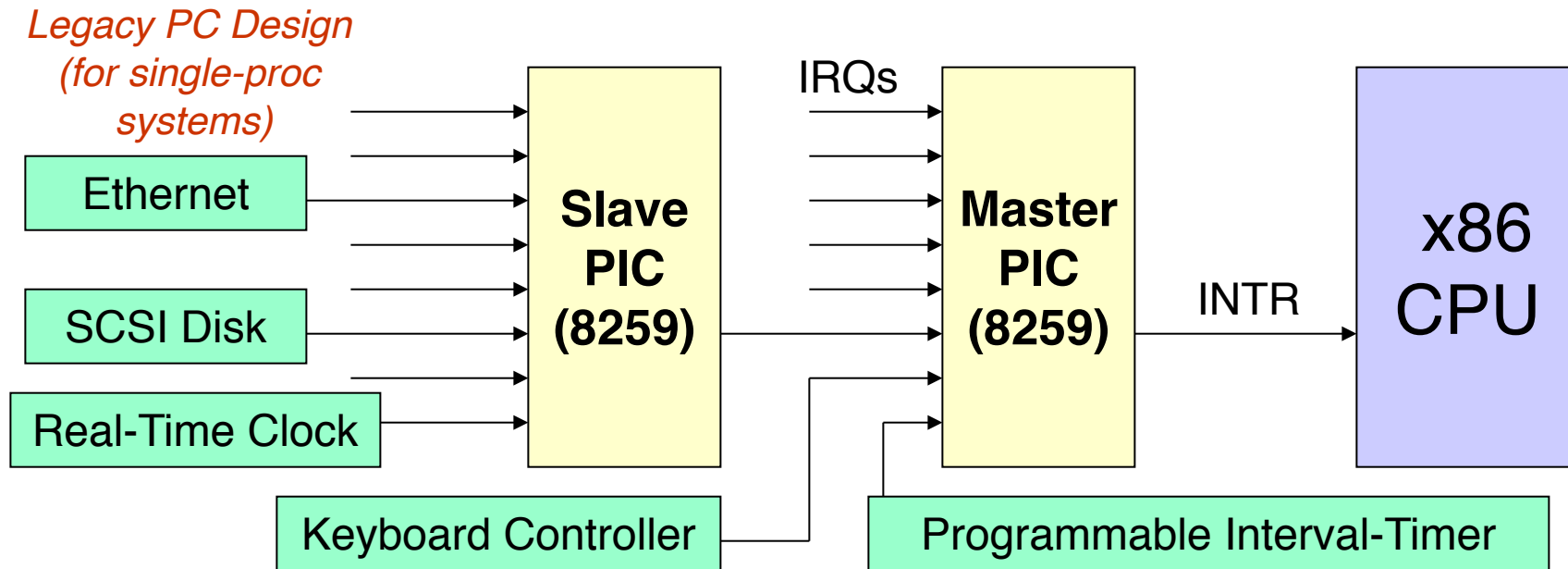
Possible Solution: Polling

- **CPU periodically checks if each device needs service**
 - ✗ takes CPU time when there are no events pending
 - ✗ can reduce checking overhead but at the expense of longer response time
 - ✓ can be efficient if events arrive rapidly

“Polling is like picking up your phone every few seconds to see if you have a call...”

“Interrupts are like waiting for the phone to ring.”

Interrupt Hardware

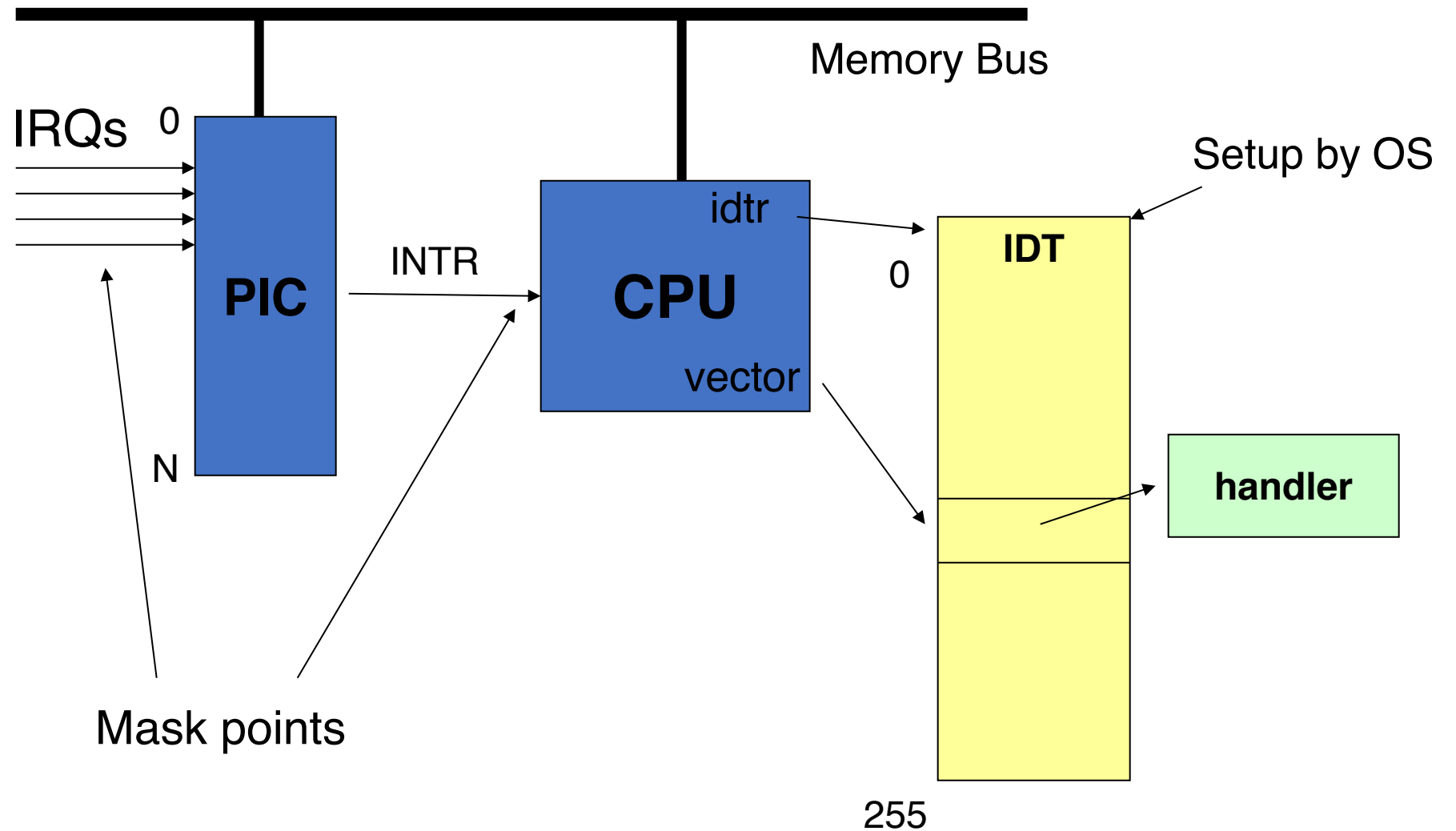


- I/O devices have (unique or shared) *Interrupt Request Lines (IRQs)*
- IRQs are mapped by special hardware to *interrupt vectors*, and passed to the CPU
- This hardware is called a *Programmable Interrupt Controller (PIC)*

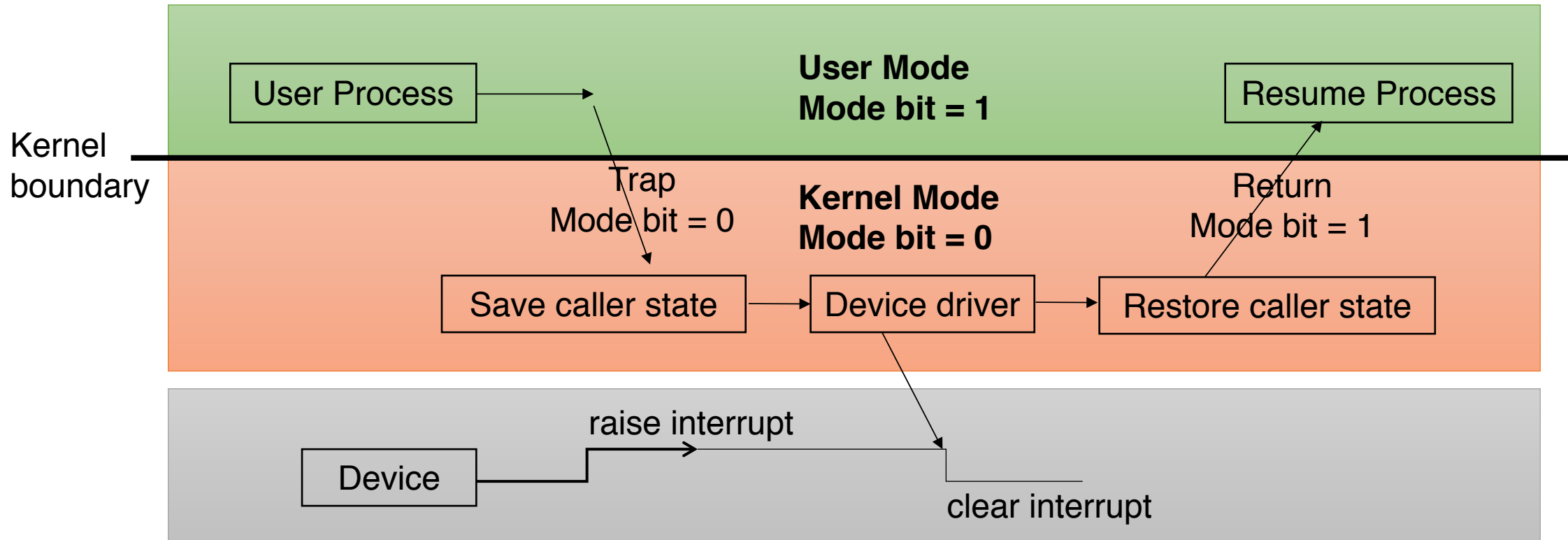
The “Interrupt Controller”

- **PIC: Programmable Interrupt Controller (8259A)**
 - Responsible for telling the CPU **when and which** device wishes to ‘interrupt’
 - Has 16 wires to devices (IRQ0 – IRQ15)
- **PIC translates IRQs to *vector* number**
 - Vector number is signaled over INTR line
 - **In Pintos:**
 - IRQ0...15 are delivered to interrupt vectors 32...47 (`src/threads/interrupt.c`)
- **Interrupts can have varying priorities**
 - PIC also needs to prioritize multiple requests
- **Possible to “mask” (disable) interrupts at PIC or CPU**

Putting It All Together



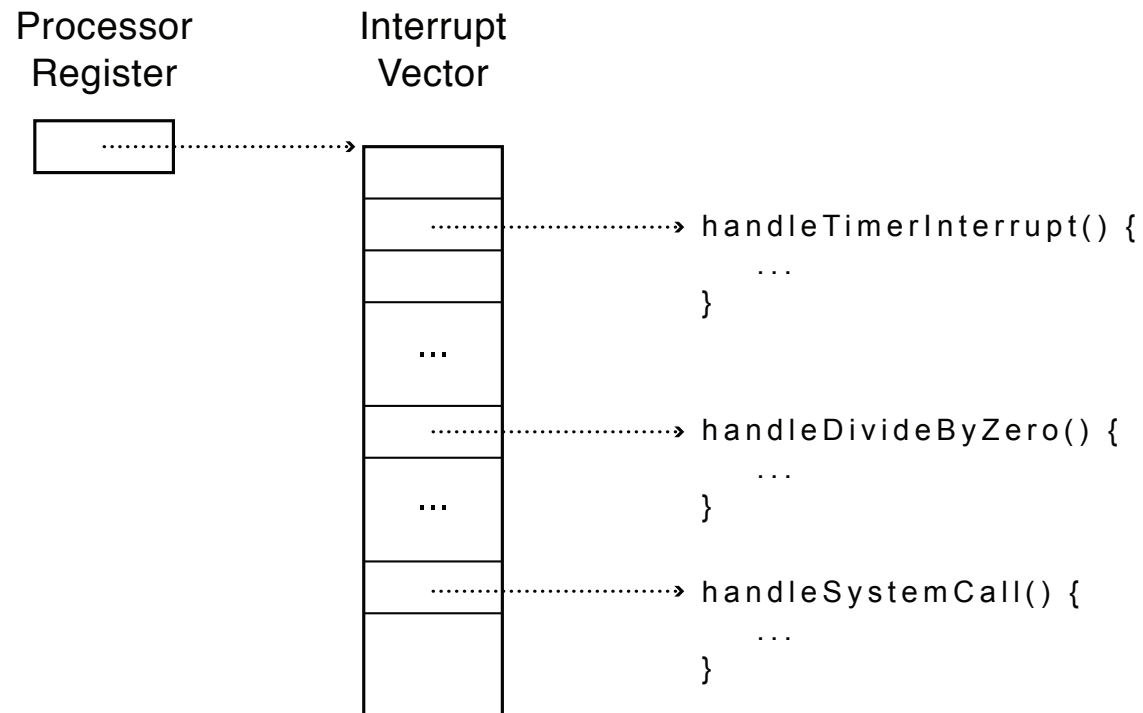
An Interrupt Illustrated



© 2004-2007 Ed Lazowska, Hank Levy, Andrea and Remzi Arpaci-Dussea, Michael Swift

Interrupt Vector Table (IVT)

- **A data structure to associate interrupt requests with handlers**
 - each entry is called an interrupt vector (specifies the address of the handler)
 - architecture-specific implementation



Interrupt Vector Table (IVT)

- **A data structure to associate interrupt requests with handlers**
 - each entry is called an interrupt vector (specifies the address of the handler)
 - architecture-specific implementation
- **In x86 called Interrupt Descriptor Table (IDT)**
 - supports 256 interrupts, so the IDT contains 256 entries
 - each entry specifies the address of the handler plus some flags
 - programmed by the OS
 - In Pintos: `make_intr_gate (src/threads/interrupt.c)`

Event: Interrupt vs. Exceptions

- Two kinds of events, interrupts and exceptions
- Interrupts are caused by an external event (asynchronous)
 - Device finishes I/O, timer expires, etc.
- **Exceptions are caused by executing instructions (synchronous)**
 - CPU requires software intervention to handle a fault or trap
 - e.g., the x86 INT instruction
 - e.g., a page fault, write to a read-only page
 - a deliberate exception is a “trap”, unexpected is a “fault”

I/O Control

- **I/O issues**
 - Initiating an I/O
 - Completing an I/O
- **Initiating an I/O**
 - Special instructions
 - Memory-mapped I/O
 - Device registers mapped into address space
 - Writing to address sends data to I/O device

I/O Completion

- **Interrupts are the basis for asynchronous I/O**
 - OS initiates I/O
 - Device operates independently of rest of machine
 - Device sends an interrupt signal to CPU when done
 - OS maintains a vector table containing a list of addresses of kernel routines to handle various events
 - CPU looks up kernel address indexed by interrupt number, context switches to routine

I/O Example

- 1. Ethernet receives packet, writes packet into memory**
- 2. Ethernet signals an interrupt**
- 3. CPU stops current operation, switches to kernel mode, saves machine state (PC, mode, etc.) on kernel stack**
- 4. CPU reads address from vector table indexed by interrupt number, branches to address (Ethernet device driver)**
- 5. Ethernet device driver processes packet (reads descriptors to find packet in memory)**
- 6. Upon completion, restores saved state from stack**

Timer

- **The timer is critical for an operating system**
- **It is the fallback mechanism for OS to reclaim control over the machine**
 - Timer is set to generate an interrupt after a period of time
 - Setting timer is a privileged instruction
 - When timer expires, generates an interrupt
 - Handled by kernel, which controls resumption context
 - Basis for OS [scheduler](#) (*more later...*)
- **Prevents infinite loops**
 - OS can always regain control from erroneous or malicious programs that try to hog CPU
- **Also used for time-based functions (e.g., *sleep()*)**

Deliberate Exception: Trap

- **A trap is an intentional software-generated exception**
 - the main mechanism for programs to interact with the OS
 - On x86, programs use the `int` instruction to cause a trap
 - On ARM, `svc` instruction
- **Handler for trap is defined in interrupt vector table**
 - Kernel chooses one vector for representing system call trap
 - e.g., `int $0x80` is used to in Linux to make system calls
 - Pintos uses `int $0x30` for system call trap

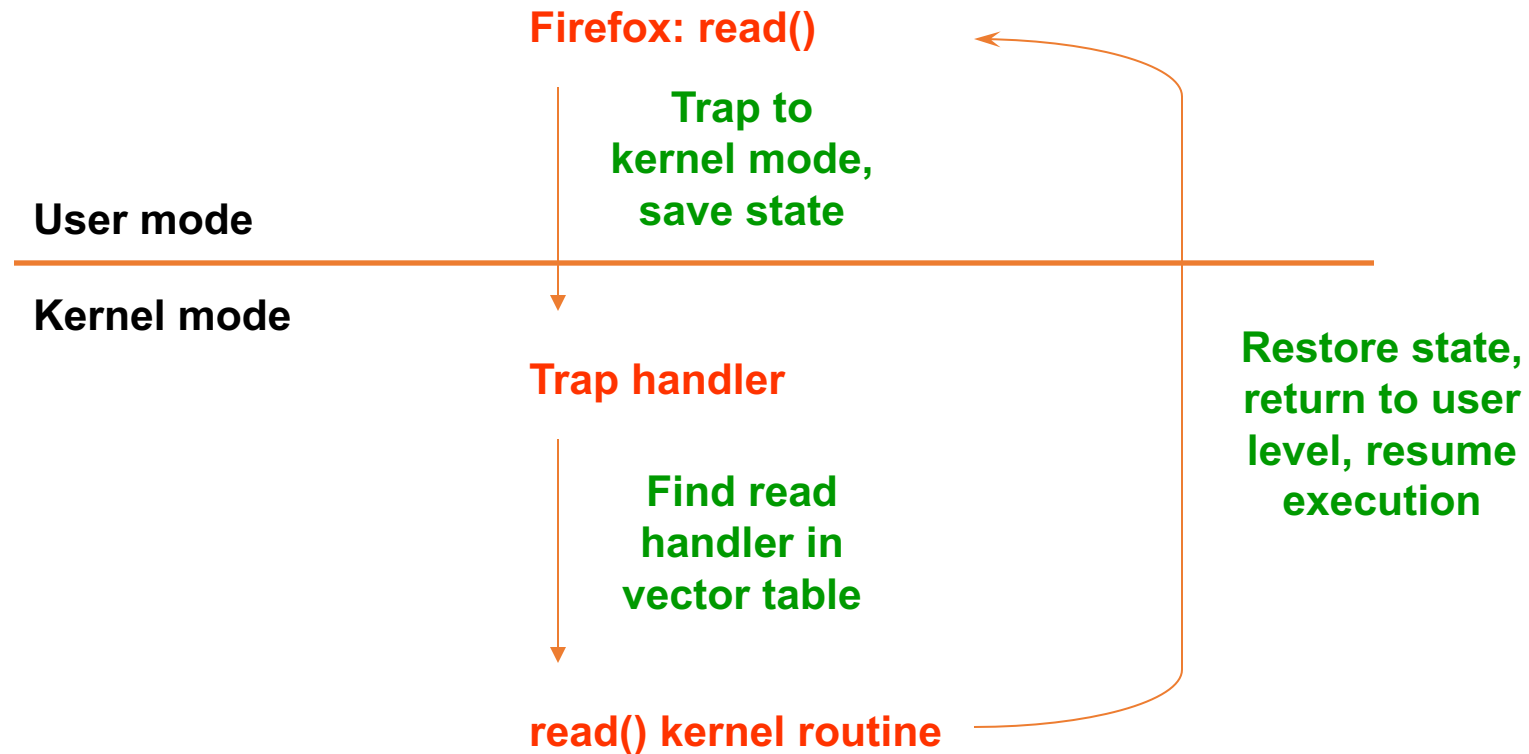
System Call Trap

- **For a user program to “call” OS service**
 - Known as **crossing the protection boundary**, or **protected control transfer**
- **The system call instruction**
 - Causes an exception, which vectors to a kernel handler
 - Passes a parameter determining the system routine to call

```
movl $20, %eax # Get PID of current process
int $0x80 # Invoke system call!
# Now %eax holds the PID of the current process
```

- Saves caller state (PC, regs, mode) so it can be restored
 - Returning from system call restores this state
- **Requires architectural support to:**
 - Restore saved state, reset mode, resume execution

System Call



LINUX System Call Quick Reference

Jialong He
Jialong_he@bigfoot.com
http://www.bigfoot.com/~jialong_he

Introduction

System call is the services provided by Linux kernel. In C programming, it often uses functions defined in **libc** which provides a wrapper for many system calls. Manual page section 2 provides more information about system calls. To get an overview, use "man 2 intro" in a command shell.

It is also possible to invoke **syscall()** function directly. Each system call has a function number defined in **<syscall.h>** or **<unistd.h>**. Internally, system call is invoked by software interrupt 0x80 to transfer control to the kernel. System call table is defined in Linux kernel source file "**arch/i386/kernel/entry.S**".

System Call Example

```
#include <syscall.h>
#include <unistd.h>
#include <stdio.h>
#include <sys/types.h>

int main(void) {

    long ID1, ID2;
    /*-----*/
    /* direct system call */
    /* SYS_getpid (func no. is 20) */
    /*-----*/
    ID1 = syscall(SYS_getpid);
    printf ("syscall(SYS_getpid)=%ld\n", ID1);

    /*-----*/
    /* "libc" wrapped system call */
    /* SYS_getpid (Func No. is 20) */
    /*-----*/
    ID2 = getpid();
    printf ("getpid()=%ld\n", ID2);

    return (0);
}
```

System Call Quick Reference

No	Func Name	Description	Source
1	exit	terminate the current process	<i>kernel/exit.c</i>
2	fork	create a child process	<i>arch/i386/kernel/process.c</i>
3	read	read from a file descriptor	<i>fs/read_write.c</i>
4	write	write to a file descriptor	<i>fs/read_write.c</i>
5	open	open a file or device	<i>fs/open.c</i>
6	close	close a file descriptor	<i>fs/open.c</i>
7	waitpid	wait for process termination	<i>kernel/exit.c</i>

8	creat	create a file or device ("man 2 open" for information)	<i>fs/open.c</i>
9	link	make a new name for a file	<i>fs/namei.c</i>
10	unlink	delete a name and possibly the file it refers to	<i>fs/namei.c</i>
11	execve	execute program	<i>arch/i386/kernel/process.c</i>
12	chdir	change working directory	<i>fs/open.c</i>
13	time	get time in seconds	<i>kernel/time.c</i>
14	mknod	create a special or ordinary file	<i>fs/namei.c</i>
15	chmod	change permissions of a file	<i>fs/open.c</i>
16	lchown	change ownership of a file	<i>fs/open.c</i>
18	stat	get file status	<i>fs/stat.c</i>
19	lseek	reposition read/write file offset	<i>fs/read_write.c</i>
20	getpid	get process identification	<i>kernel/sched.c</i>
21	mount	mount filesystems	<i>fs/super.c</i>
22	umount	unmount filesystems	<i>fs/super.c</i>
23	setuid	set real user ID	<i>kernel/sys.c</i>
24	getuid	get real user ID	<i>kernel/sched.c</i>
25	stime	set system time and date	<i>kernel/time.c</i>
26	ptrace	allows a parent process to control the execution of a child process	<i>arch/i386/kernel/ptrace.c</i>
27	alarm	set an alarm clock for delivery of a signal	<i>kernel/sched.c</i>
28	fstat	get file status	<i>fs/stat.c</i>
29	pause	suspend process until signal	<i>arch/i386/kernel/sys_i386.c</i>
30	utime	set file access and modification times	<i>fs/open.c</i>
33	access	check user's permissions for a file	<i>fs/open.c</i>
34	nice	change process priority	<i>kernel/sched.c</i>
36	sync	update the super block	<i>fs/buffer.c</i>
37	kill	send signal to a process	<i>kernel/signal.c</i>
38	rename	change the name or location of a file	<i>fs/namei.c</i>
39	mkdir	create a directory	<i>fs/namei.c</i>
40	rmdir	remove a directory	<i>fs/namei.c</i>
41	dup	duplicate an open file descriptor	<i>fs/fcntl.c</i>
42	pipe	create an interprocess channel	<i>arch/i386/kernel/sys_i386.c</i>
43	times	get process times	<i>kernel/sys.c</i>
45	brk	change the amount of space allocated for the calling process's data segment	<i>mm/mmap.c</i>
46	setgid	set real group ID	<i>kernel/sys.c</i>
47	getgid	get real group ID	<i>kernel/sched.c</i>
48	sys_signal	ANSI C signal handling	<i>kernel/signal.c</i>
49	geteuid	get effective user ID	<i>kernel/sched.c</i>
50	getegid	get effective group ID	<i>kernel/sched.c</i>

System Call Questions

- What would happen if the kernel did not save state?
- What if the *kernel* executes a system call?
- What if a user program returns from a system call?
- How to reference kernel objects as arguments or results to/from syscalls?
 - A naming issue
 - Use integer object handles or descriptors
 - E.g., Unix file descriptors, Windows HANDLEs
 - Only meaningful as parameters to other system calls
 - Also called capabilities (more later when we do protection)
 - Why not use kernel addresses to name kernel objects?

Unexpected Exception: Faults

- **Hardware detects and reports “exceptional” conditions**
 - Page fault, unaligned access, divide by zero
- **Upon exception, hardware “faults” (verb)**
 - Must save state (PC, regs, mode, etc.) so that the faulting process can be restarted
- **Modern OSes use VM faults for many functions**
 - Debugging, end-of-stack, garbage collection, copy-on-write
- **Fault exceptions are a performance optimization**
 - Could detect faults by inserting extra instructions into code (at a significant performance penalty)

Handling Faults

- **Some faults are handled by “fixing”...**
 - “Fix” the exceptional condition and return to the faulting context
 - Page faults cause the OS to place the missing page into memory
 - Fault handler resets PC of faulting context to **re-execute** instruction that caused the page fault
- **Some faults are handled by notifying the process**
 - Fault handler changes the saved context to transfer control to a user-mode handler on return from fault
 - Handler must be registered with OS
 - Unix **signals** or Win **user-mode Async Procedure Calls (APCs)**
 - SIGALRM, SIGHUP, SIGTERM, SIGSEGV, etc.

Handling Faults (2)

- **Kernel may handle unrecoverable faults by killing the process**
 - Program fault with no registered handler
 - Halt process, write process state to file, destroy process
 - In Unix, the default action for many signals (e.g., SIGSEGV)
- **What about faults in the kernel?**
 - Dereference NULL, divide by zero, undefined instruction
 - These faults considered fatal, operating system crashes
 - **Unix panic**, **Windows “Blue screen of death”**
 - Kernel is halted, state dumped to a core file, machine locked up

Types of Arch Support

- **Manipulating privileged machine state**
 - Protected instructions
 - Manipulate device registers, TLB entries, etc.
- **Generating and handling “events”**
 - Interrupts, exceptions, system calls, etc.
 - Respond to external events
 - CPU requires software intervention to handle fault or trap
- **Mechanisms to handle concurrency**
 - Interrupts, atomic instructions

Synchronization

- **Interrupts cause difficult problems**
 - An interrupt can occur at any time
 - A handler can execute that interferes with code that was interrupted
- **OS must be able to synchronize concurrent execution**
- **Need to guarantee that short instruction sequences execute atomically**
 - Disable interrupts – turn off interrupts before sequence, execute sequence, turn interrupts back on
 - Special atomic instructions – read/modify/write a memory address, test and conditionally set a bit based upon previous value
 - XCHG instruction on x86

Summary

- **Protection**

- User/kernel modes
- Protected instructions

- **Interrupts**

- Timer, I/O

	Unexpected	Deliberate
Exceptions (sync)	fault	syscall trap
Interrupts (async)	interrupt	software interrupt

- **System calls**

- Used by user-level processes to access OS functions
- Access what is “in” the OS

- **Exceptions**

- Unexpected event during execution (e.g., divide by zero)

Next Time...

- **Read Chapters 4-6 (Processes)**
- **Homework #1**
- **Lab 0**