

Stony Brook University CSE 506: Operating Systems

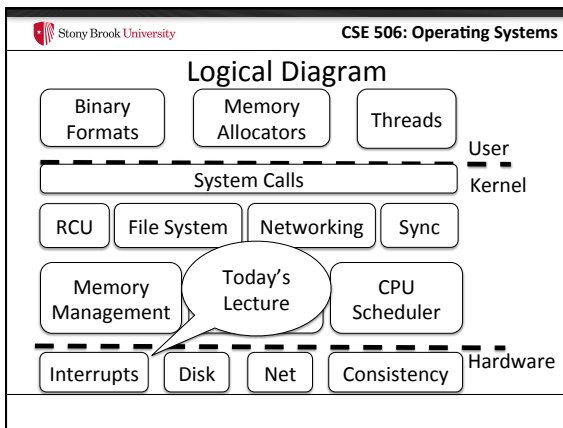
## Interrupts and System Calls

Don Porter  
CSE 506

Stony Brook University CSE 506: Operating Systems

## Housekeeping

- Welcome Tas Amogh Akshintala and Yizheng Jiao (1/2 time) – Office Hours posted
- Next Thursday’s class has a reading assignment
- Lab 1 due Monday 9/8
- All students should have VMs and private repos soon
  - Email Don if you don’t have one by tonight
  - Unless you just got in the class since Wed night (sigh)



Stony Brook University CSE 506: Operating Systems

### Background: Control Flow

```

pc // x = 2, y = void printf(va_args)
true {
  if (y) { //...
    2 /= x;
    printf(x);
  } //...

```

Regular control flow: branches and calls (logically follows source code)

Stony Brook University CSE 506: Operating Systems

### Background: Control Flow

```

pc // x = 0 void handle_divzero()
true {
  if (y) {
    2 /= x;
    printf(x);
  } //...

```

Divide by zero!  
Program can't make progress!

Irregular control flow: exceptions, system calls, etc.

Stony Brook University CSE 506: Operating Systems

### Lecture goal

- Understand the hardware tools available for irregular control flow.
  - I.e., things other than a branch in a running program
- Building blocks for context switching, device management, etc.

Stony Brook University CSE 506: Operating Systems

### Two types of interrupts

- Synchronous: will happen every time an instruction executes (with a given program state)
  - Divide by zero
  - System call
  - Bad pointer dereference
- Asynchronous: caused by an external event
  - Usually device I/O
  - Timer ticks (well, clocks can be considered a device)

Stony Brook University CSE 506: Operating Systems

### Intel nomenclature

- Interrupt – only refers to asynchronous interrupts
- Exception – synchronous control transfer

• Note: from the programmer’s perspective, these are handled with the same abstractions

Stony Brook University CSE 506: Operating Systems

### Lecture outline

- Overview
- How interrupts work in hardware
- How interrupt handlers work in software
- How system calls work
- New system call hardware on x86

Stony Brook University CSE 506: Operating Systems

### Interrupt overview

- Each interrupt or exception includes a number indicating its type
- E.g., 14 is a page fault, 3 is a debug breakpoint
- This number is the index into an interrupt table

Stony Brook University CSE 506: Operating Systems

### x86 interrupt table

The diagram shows a horizontal bar representing the x86 interrupt table, divided into segments. The first segment, from index 0 to 31, is labeled 'Reserved for the CPU'. The second segment, from index 32 to 47, is labeled 'Device IRQs'. The third segment, from index 48 to 255, is labeled 'Software Configurable'. Within the 'Software Configurable' segment, two callout boxes point to specific indices: '48 = JOS System Call' and '128 = Linux System Call'. Ellipses (...) are used to indicate that the table continues between 31 and 47, and between 47 and 255.

Stony Brook University CSE 506: Operating Systems

### x86 interrupt overview

- Each type of interrupt is assigned an index from 0—255.
- 0—31 are for processor interrupts; generally fixed by Intel
  - E.g., 14 is always for page faults
- 32—255 are software configured
  - 32—47 are for device interrupts (IRQs) in JOS
    - Most device’s IRQ line can be configured
    - Look up APICs for more info (Ch 4 of Bovet and Cesati)
  - 0x80 issues system call in Linux (more on this later)

Stony Brook University CSE 506: Operating Systems

### Software interrupts

- The `int <num>` instruction allows software to raise an interrupt
  - `0x80` is just a Linux convention. JOS uses `0x30`.
- There are a lot of spare indices
  - You could have multiple system call tables for different purposes or types of processes!
    - Windows does: one for the kernel and one for win32k

Stony Brook University CSE 506: Operating Systems

### Software interrupts, cont

- OS sets ring level required to raise an interrupt
  - Generally, user programs can't issue an `int 14` (page fault) manually
  - An unauthorized `int` instruction causes a general protection fault
    - Interrupt 13

Stony Brook University CSE 506: Operating Systems

### What happens (generally):

- Control jumps to the kernel
  - At a prescribed address (the interrupt handler)
- The register state of the program is dumped on the kernel's stack
  - Sometimes, extra info is loaded into CPU registers
  - E.g., page faults store the address that caused the fault in the `cr2` register
- Kernel code runs and handles the interrupt
- When handler completes, resume program (see `iret` instr.)

Stony Brook University CSE 506: Operating Systems

### How it works (HW)

- How does HW know what to execute?
- Where does the HW dump the registers; what does it use as the interrupt handler's stack?

Stony Brook University CSE 506: Operating Systems

### How is this configured?

- Kernel creates an array of Interrupt descriptors in memory, called Interrupt Descriptor Table, or IDT
  - Can be anywhere in memory
  - Pointed to by special register (`idt_r`)
    - c.f., segment registers and `gdt_r` and `ldt_r`
- Entry 0 configures interrupt 0, and so on

Stony Brook University CSE 506: Operating Systems

### x86 interrupt table

The diagram illustrates the x86 interrupt table as a sequence of 256 entries. The first entry at index 0 is pointed to by the `idt_r` register. A callout box indicates that the linear address of the interrupt table starts at index 0. The entries are indexed from 0 to 255, with ellipses indicating intermediate entries between 31 and 47, and between 47 and 255.

Stony Brook University CSE 506: Operating Systems

### x86 interrupt table

```

Code Segment: Kernel Code
Segment Offset: &page_fault_handler //linear addr
Ring: 0 // kernel
Present: 1
Gate Type: Exception

```

Stony Brook University CSE 506: Operating Systems

### Interrupt Descriptor

- Code segment selector
  - Almost always the same (kernel code segment)
  - Recall, this was designed before paging on x86!
- Segment offset of the code to run
  - Kernel segment is “flat”, so this is just the linear address
- Privilege Level (ring)
  - Ring that can raise this interrupt with an `int` instruction
- Present bit – disable unused interrupts
- Gate type (interrupt or trap/exception) – more in a bit

Stony Brook University CSE 506: Operating Systems

### x86 interrupt table

```

Code Segment: Kernel Code
Segment Offset: &breakpoint_handler //linear addr
Ring: 3 // user
Present: 1
Gate Type: Exception

```

Stony Brook University CSE 506: Operating Systems

### Interrupt Descriptors, ctd.

- In-memory layout is a bit confusing
  - Like a lot of the x86 architecture, many interfaces were later deprecated
- Worth comparing Ch 9.5 of the i386 manual with `inc/mmu.h` in the JOS source code

Stony Brook University CSE 506: Operating Systems

### How it works (HW)

- How does HW know what to execute?
  - Interrupt descriptor table specifies what code to run
    - And at what privilege (via code segment)
  - This can be set up once during boot for the whole system
- Where does the HW dump the registers; what does it use as the interrupt handler’s stack?
  - Specified in the Task State Segment

Stony Brook University CSE 506: Operating Systems

### Task State Segment (TSS)

- Another segment, just like the code and data segment
  - A descriptor created in the GDT (cannot be in LDT)
  - Selected by special task register (`tr`)
  - Unlike others, has a hardware-specified layout
- Lots of fields for rarely-used features
- Two features we care about in a modern OS:
  - 1) Location of kernel stack (fields `ss0/esp0`)
  - 2) I/O Port privileges (more in a later lecture)

Stony Brook University CSE 506: Operating Systems

### TSS, cont.

- Simple model: specify a TSS for each process
  - Note: Only 2<sup>13</sup> entries in the GDT
- Optimization (JOS):
  - Our kernel is pretty simple (uniprocessor only)
  - Why not just share one TSS and kernel stack per-process?
- Linux generalization:
  - One TSS per CPU
  - Modify TSS fields as part of context switching

Stony Brook University CSE 506: Operating Systems

### Summary

- Most interrupt handling hardware state set during boot
- Each interrupt has an IDT entry specifying:
  - What code to execute, privilege level to raise the interrupt
- Stack to use specified in the TSS

Stony Brook University CSE 506: Operating Systems

### Comment

- Again, segmentation rears its head
- You can't program OS-level code on x86 without getting your hands dirty with it
- Helps to know which features are important when reading the manuals

Stony Brook University CSE 506: Operating Systems

### Lecture outline

- Overview
- How interrupts work in hardware
- **How interrupt handlers work in software**
- How system calls work
- New system call hardware on x86

Stony Brook University CSE 506: Operating Systems

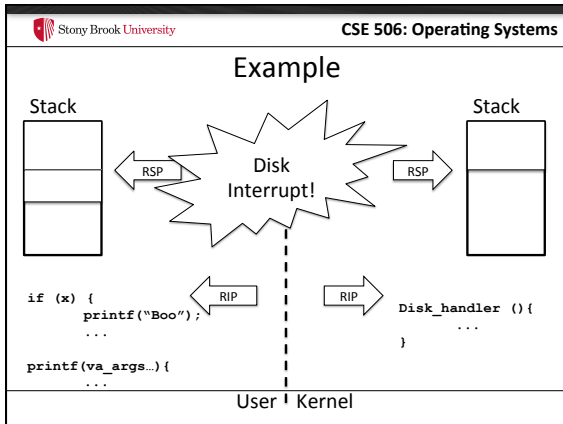
### High-level goal

- Respond to some event, return control to the appropriate process
- What to do on:
  - Network packet arrives
  - Disk read completion
  - Divide by zero
  - System call

Stony Brook University CSE 506: Operating Systems

### Interrupt Handlers

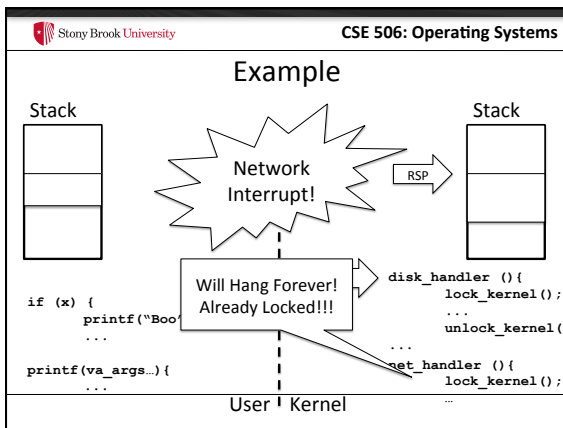
- Just plain old kernel code



Stony Brook University CSE 506: Operating Systems

### Complication:

- What happens if I'm in an interrupt handler, and another interrupt comes in?
  - Note: kernel stack only changes on privilege level change
  - Nested interrupts just push the next frame on the stack
- What could go wrong?
  - Violate code invariants
  - Deadlock
  - Exhaust the stack (if too many fire at once)



Stony Brook University CSE 506: Operating Systems

### Bottom Line:

- Interrupt service routines must be reentrant or synchronize
- Period.

Stony Brook University CSE 506: Operating Systems

### Hardware interrupt sync.

- While a CPU is servicing an interrupt on a given IRQ line, the same IRQ won't raise another interrupt until the routine completes
  - Bottom-line: device interrupt handler doesn't have to worry about being interrupted by itself
- A *different* device can interrupt the handler
  - Problematic if they share data structures
  - Like a list of free physical pages...
  - What if both try to grab a lock for the free list?

Stony Brook University CSE 506: Operating Systems

### Disabling interrupts

- An x86 CPU can disable I/O interrupts
  - Clear bit 9 of the EFLAGS register (IF Flag)
  - `cld` and `sti` instructions clear and set this flag
- Before touching a shared data structure (or grabbing a lock), an interrupt handler should disable I/O interrupts

Stony Brook University CSE 506: Operating Systems

## Gate types

- Recall: an IDT entry can be an interrupt or an exception gate
- Difference?
  - An interrupt gate automatically disables all other interrupts (i.e., clears and sets IF on enter/exit)
  - An exception gate doesn't
- This is just a programmer convenience: you could do the same thing in software

Stony Brook University CSE 506: Operating Systems

## Exceptions

- You can't mask exceptions
  - Why not?
    - Can't make progress after a divide-by-zero
  - Double and Triple faults detect faults in the kernel
- Do exception handlers need to be reentrant?
  - Not if your kernel has no bugs (or system calls in itself)
  - In certain cases, Linux allows nested page faults
    - E.g., to detect errors copying user-provided buffers

Stony Brook University CSE 506: Operating Systems

## Summary

- Interrupt handlers need to synchronize, both with locks (multi-processor) and by disabling interrupts (same CPU)
- Exception handlers can't be masked
  - Nested exceptions generally avoided

Stony Brook University CSE 506: Operating Systems

## Lecture outline

- Overview
- How interrupts work in hardware
- How interrupt handlers work in software
- **How system calls work**
- New system call hardware on x86

Stony Brook University CSE 506: Operating Systems

## System call "interrupt"

- Originally, system calls issued using `int` instruction
- Dispatch routine was just an interrupt handler
- Like interrupts, system calls are arranged in a table
  - See `arch/x86/kernel/syscall_table*.S` in Linux source
- Program selects the one it wants by placing index in `eax` register
  - Arguments go in the other registers by calling convention
  - Return value goes in `eax`

Stony Brook University CSE 506: Operating Systems

## Lecture outline

- Overview
- How interrupts work in hardware
- How interrupt handlers work in software
- How system calls work
- **New system call hardware on x86**

Stony Brook University CSE 506: Operating Systems

### Around P4 era...

- Processors got very deeply pipelined
  - Pipeline stalls/flushes became very expensive
  - Cache misses can cause pipeline stalls
- System calls took twice as long from P3 to P4
  - Why?
  - IDT entry may not be in the cache
  - Different permissions constrain instruction reordering

Stony Brook University CSE 506: Operating Systems

### Idea

- What if we cache the IDT entry for a system call in a special CPU register?
  - No more cache misses for the IDT!
  - Maybe we can also do more optimizations
- Assumption: system calls are frequent enough to be worth the transistor budget to implement this
  - What else could you do with extra transistors that helps performance?

Stony Brook University CSE 506: Operating Systems

### AMD: syscall/sysret

- These instructions use MSRs (machine specific registers) to store:
  - syscall entry point and code segment
  - Kernel stack
- A drop-in replacement for `int 0x80`
- Everyone loved it and adopted it wholesale
  - Even Intel!

Stony Brook University CSE 506: Operating Systems

### Aftermath

- Getpid() on my desktop machine (recent AMD 6-core):
  - int 80: 371 cycles
  - syscall: 231 cycles
- So system calls are definitely faster as a result!

Stony Brook University CSE 506: Operating Systems

### In JOS

- You will use the `int` instruction to implement system calls
- There is a challenge problem in lab 3 (i.e., extra credit) to use `sysenter/sysexit`
  - Note that there are some more details about register saving to deal with
  - syscall/sysret is a bit too trivial for extra credit
    - But still cool if you get it working!

Stony Brook University CSE 506: Operating Systems

### Summary

- Interrupt handlers are specified in the IDT
- Understand when nested interrupts can happen
  - And how to prevent them when unsafe
- Understand optimized system call instructions
  - Be able to explain `vdso`, `syscall` vs. `int 80`