# ECE391 Computer System Engineering Lecture 12

Dr. Zbigniew Kalbarczyk
University of Illinois at Urbana- Champaign

Fall 2018

### Lecture Topics

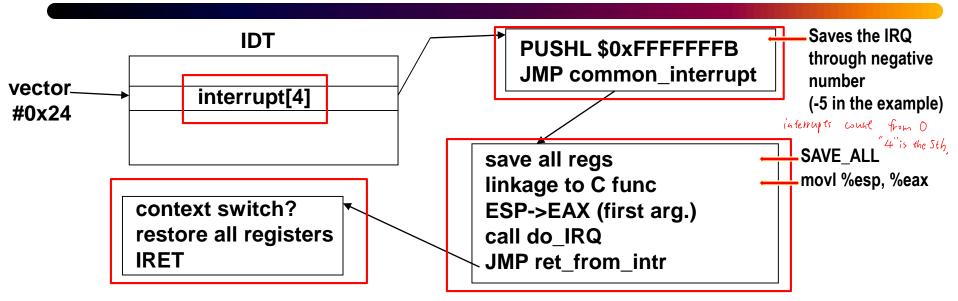
- Linux interrupt system
  - data structures
  - handler installation & removal
  - invocation
  - execution
  - soft interrupts (tasklets) in Linux

### **Aministrivia**

- MP2 Checkpoint 1
  - Due by 6:00pm Monday, October 8

- MP3 Teams
  - Due by end of Friday, October 12

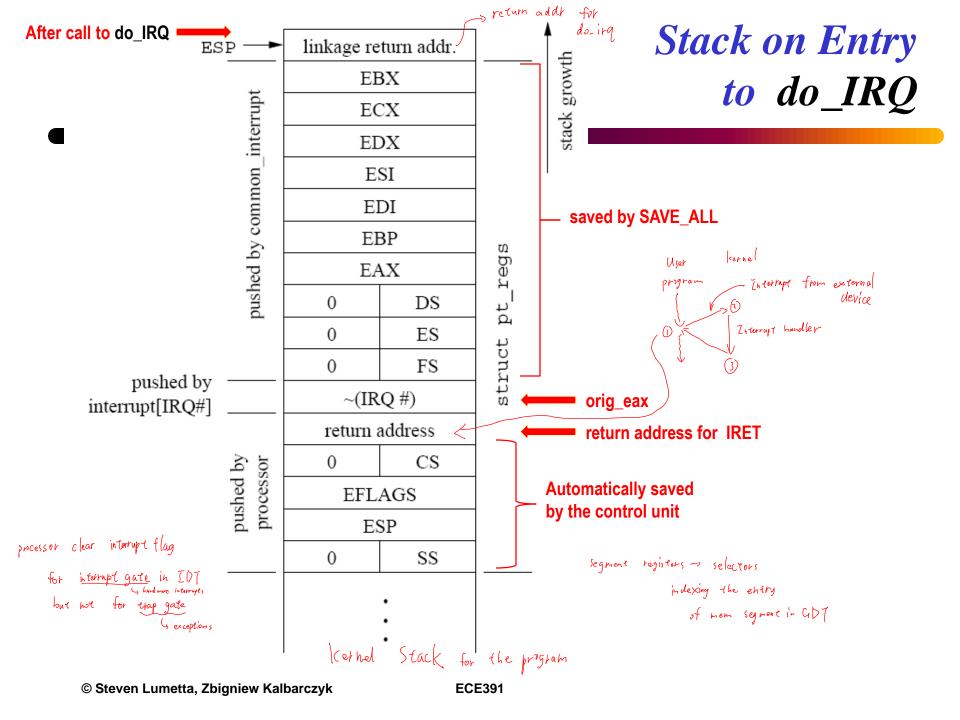
### Interrupt Invocation



- 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

	0x00	division error	
	0x02	NMI (	
	0x02	• • • • • • • • • • • • • • • • • • • •	
0x00-0x1F	0x03	breakpoint (used by KGDB) overflow	
0x00=0x1F	0x04	overnow	
	:		
defined	0x0B	segment not present	
by Intel	0x0C	stack segment fault	
	0x0D	general protection fault	
	0x0E	page fault	
	:		
	0x20	IRQ0 — timer chip	
	0x21		
0x20-0x27	0x22	IRQ2 — (cascade to slave)	
	ı	IRQ3	
master	0x24	IRQ4 — serial port (KGDB)	
8259 PIC	0x25	IRQ5	
	0x26	IRQ6	example
	1	IRQ7	of
	1	IRQ8 — real time clock	possible
	1	IRQ9	settings
0x28-0x2F	I	IRQ10	
		IRQ11 — eth0 (network)	
slave	1	IRQ12 — PS/2 mouse	
8259 PIC		IRQ13	
	1	IRQ14 — ide0 (hard drive)	
	0x2F	IRQ15	
0x30-0x7F	:	APIC vectors available to device drivers	
0x80	0x80	system call vector (INT 0x80)	
0x81-0xEE	:	more APIC vectors available to device drivers	
0xEF	0xEF	local APIC timer	
0xF0-0xFF	:	symmetric multiprocessor (SMP) communication vectors	

## Interrupt Descriptor Table



# Linux do\_IRQ arch/i386/kernel/irq.c

```
fastcall unsigned int
                                            fastcall convention used to pass arguments in registers;
                                            reduces the number of memory accesses required for the call
do IRQ (struct pt regs* regs)
    struct pt regs* old regs;
    /* high bit used in ret from code */
    struct irq desc* desc = irq desc + irq;
    if ((unsigned)irq >= NR IRQS) {
        printk (KERN EMERG "%s: cannot handle IRQ %d\n", FUNCTION , irq);
        BUG ();
    old regs = set irq regs (regs); Record registers at time of interrupt
    irq enter ();
                                         Increments counter of nested interrupt handlers
    /* for 8259A interrupts, handle irg is set to handle level irg */
    desc->handle irq (irq, desc); | call interrupt flow handler handle_level_irq for all 8259A interrupts
                                       Exit irg context and process softirgs if needed
    irq exit ();
    set irq regs (old regs);
    return 1;
```

### Interrupt Invocation (cont.)

Signature for do\_IRQ

```
fastcall unsigned int do_IRQ
  (struct pt_regs* regs);
```

- fastcall macro in Linux tells gcc to pass args in EAX, EDX, ECX
- EAX points to saved registers (regs argument)
- saved registers already on stack
- note that processor clears IF when it takes an interrupt
- EFLAGS stored by processor include original IF value

### Comments on do\_irq

Value pushed by irq-specific code re-converted to find irq #

Start with a sanity check

- set\_irq\_regs calls record registers at time of interrupt
  - per-CPU storage

### Comments on do\_irq (cont.)

- irq\_enter / irq\_exit
  - necessary for proper priority
    - e.g., second hard interrupt occurs
    - after new interrupt handled, should return to first
    - must delay processing of soft interrupts
  - irq exit processes soft interrupts when appropriate

### Comments on do\_irq (cont.)

- Central component
  - call interrupt flow handler
  - handle\_level\_irq for all 8259A interrupts
- Return value ignored (probably intended to indicate interrupt handled)

```
fastcall void
handle level irq (unsigned int irq, struct irq desc* desc)
                                                                     Linux handle_level_irq
    unsigned int
                        cpu = smp processor id ();
    struct irqaction* action;
                                                                                         kernel/irq/chip.c
    irgreturn t
                        action ret;
                                                Critical section begins
    spin lock (&desc->lock);
    mask ack irq (desc, irq);
                                                Call PIC's mask ack function
    if (desc->status & IRQ INPROGRESS)
                                                       If interrupt already in progress, do nothing
        goto out unlock;
    desc->status &= ~(IRQ REPLAY
                                       IRQ WAITING);
    kstat cpu (cpu).irqs[irq]++;
    /*
     * If its disabled or no action available
     * keep it masked and get out of here
    action = desc->action;
    if (!action | | (desc->status & IRQ DISABLED))
                                                                     Atomically decide whether to execute handlers
         desc->status |= IRQ PENDING;
         goto out unlock;
    desc->status |= IRQ INPROGRESS;
                                                   Set status for execution: in-progress and no longer pending
    desc->status &= ~IRQ PENDING;
    spin unlock (&desc->lock);
                                                 Critical section ends
    action ret = handle IRQ event (irq, action);
                                                                        Handler execution done via handle IRQ event
    if (!noirgdebug)
        note interrupt (irq, desc, action ret);
                                                  Critical section begins
    spin lock(&desc->lock);
    desc->status &= ~IRQ INPROGRESS;
    if (!(desc->status & IRQ DISABLED) && desc->chip->unmask)
                                                                               When done, remove in-progress flag
        desc->chip->unmask (irq);
                                                                               and unmask on PIC (unless disabled)
out unlock:
                                                  Critical section ends
ECE391
    spin unlock (&desc->lock);

© Steven Lumetta, Zbigniew Kalbarczyk
```

### Comments on handle\_level\_irq

- Critical section starts
  - to read descriptor status and action (handler) list
  - IF=0 at this point (set by processor when taking interrupt)
- Immediately call PIC's mask\_ack function (via a wrapper function)
- If interrupt already in progress, do nothing
  - Interrupt already re-masked and re-acknowledged on PIC
- Remove software replay and autoprobe flags
- Kernel statistics track # of interrupts seen (see /proc/interrupts)

### Comments on handle\_level\_irq (cont.)

- Atomically decide whether to execute handlers immediately
  - do so if handler defined and not disabled by software
  - if not, skip to end
    - mark as pending and end interrupt handling
    - replayed after last nested enable irq call

 Set status for execution: in-progress, and no longer pending

### Comments on handle\_level\_irq (cont.)

- Handler execution done via handle IRQ event
  - done without descriptor lock
    - allows handlers to use infrastructure
    - e.g., enable\_irq/disable\_irq, request\_irq/free\_irq
  - usually done with IF=1 (changed inside handle\_IRQ\_event)
    - done without descriptor lock
    - otherwise this code can deadlock with self

 When done, remove in-progress flag and unmask on PIC (unless disabled)

# Linux handle\_IRQ\_event kernel/irg/handle.c

```
irgreturn t
handle IRQ event (unsigned int irq, struct irqaction* action)
    irqreturn t ret;
    irgreturn t retval = IRQ NONE;
    unsigned int status = 0;
    /* call specific to ARM processor (to disambiguate timer ticks) */
    handle dynamic tick (action);
    if (!(action->flags & IRQF DISABLED))
                                                          call translates basically to STI
         local irq enable in hardirq ();
                                                          (local means "on this CPU")
                                                                 Walk through list; call each handler
    do
                                                                  (return value 1 means handled)
         ret = action->handler (irq, action->dev id);
         if (ret == IRQ HANDLED)
              status |= action->flags;
         retval |= ret;
         action = action->next;
    } while (action);
    if (status & IRQF SAMPLE RANDOM)
         add interrupt randomness (irg);
    local irq disable ();
                                                     Turn interrupts back off (IF=0 / CLI)
    return retval;
   © Steven Lumetta, Zbigniew Kalbarczyk
                                         ECE391
```

### Comments on handle IRQ event

- Set IF=1 unless the first handler asked to be executed with IF=0
  - indicated by IRQF\_DISABLED flag
  - call translates basically to STI (local means "on this CPU")
- Walk through list
  - call each handler
  - return value 1 means handled
- Generate random numbers if desired
- Turn interrupts back off (IF=0 / CLI)

### Summary on Descriptor Flags

- IRQ\_PENDING—interrupt raised by hardware; waiting to be executed
- IRQ\_INPROGRESS—some processor is executing handlers
- IRQ\_DISABLED—interrupt disabled in software; postpone execution
- IRQ\_REPLAY—software replay of previously postponed execution

# Soft Interrupts in Linux linux/interrupt.h and kernel/softirq.c

- When are soft interrupts executed?
  - after a hard interrupt completes
  - periodically by a daemon in the kernel

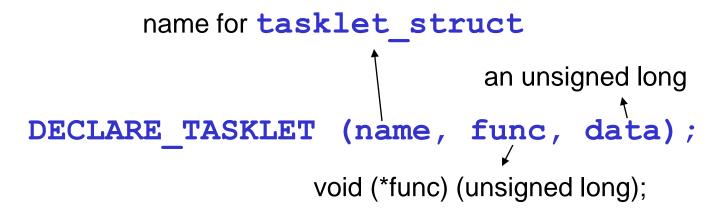
#### How?

- seven or eight prioritized types, including high and low tasklet priorities
- linked list for each tasklet priority
- run on processor on which interrupt is scheduled
- each handler atomic with respect to itself (only)

### Soft Interrupts in Linux (cont.)

linux/interrupt.h and kernel/softirq.c

#### Declaring a handler



next -	→ linked list
state	TASKLET_STATE_SCHED, TASKLET_STATE_RUN
count	# of disables
func	pointer to the tasklet function
data	integer which can be used by the tasklet function

### Tasklet Scheduling

The following two calls schedule a tasklet for execution

```
void tasklet_schedule (struct tasklet_struct* t);
void tasklet_hi_schedule (struct tasklet_struct* t);
```

- First form
  - schedules tasklet at low priority
  - on the executing processor
- Second form schedules at high priority
- Enable and disable calls analogous to hard interrupts (including nesting)

#### Tasklet Execution

- do\_softirq call
  - checks per-processor bit vector of pending priorities (high, low, etc.)
  - executes action for each priority [softirq\_vec]
  - tasklet\_action walks through linked list
    [tasklet\_hi\_action walks through high-priority list]
  - repeats up to 10 times or until no softirqs are raised

### Tasklet Execution Atomicity

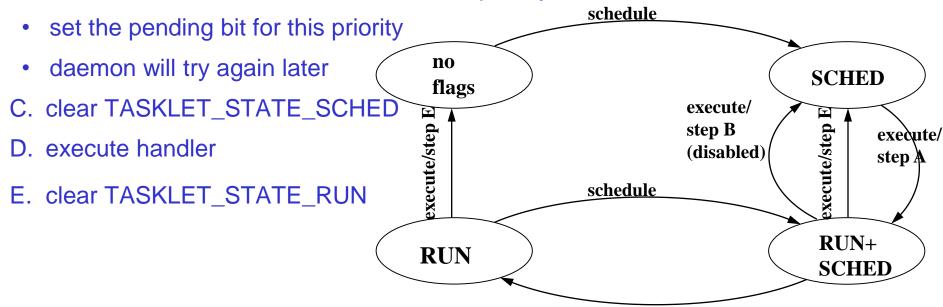
- Two bits in state changed atomically
  - TASKLET\_STATE\_SCHED tasklet scheduled for execution
  - TASKLET\_STATE\_RUN tasklet executing on some processor

- When scheduling
  - set TASKLET\_STATE\_SCHED atomically
  - if already set, schedule call does nothing

### Tasklet Execution Atomicity (cont.)

When executing, for each tasklet in linked list (at either priority)

- A. set TASKLET\_STATE\_RUN atomically (if already set, stop)
- B. check if tasklet is software disabled (count field)
  - if so, clear TASKLET\_STATE\_RUN
  - leave the tasklet in the linked list for this priority



execute/step C execution (execute/step D) occurs in lower two states (and execute/step A fails in these states)