# ECE391 Computer System Engineering Lecture 13

Dr. Zbigniew Kalbarczyk
University of Illinois at Urbana- Champaign

Spring 2021

#### Lecture Topics

- Linux interrupt system
  - data structures
  - handler installation & removal
  - invocation
  - execution

Summary of Linux interrupt system

#### **Aministrivia**

- MP2 Checkpoint 2
  - Due by 5:59pm Monday, March 15

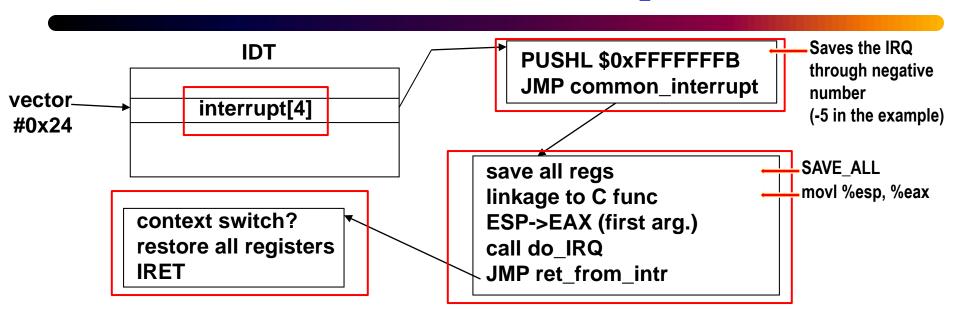
#### **Aministrivia**

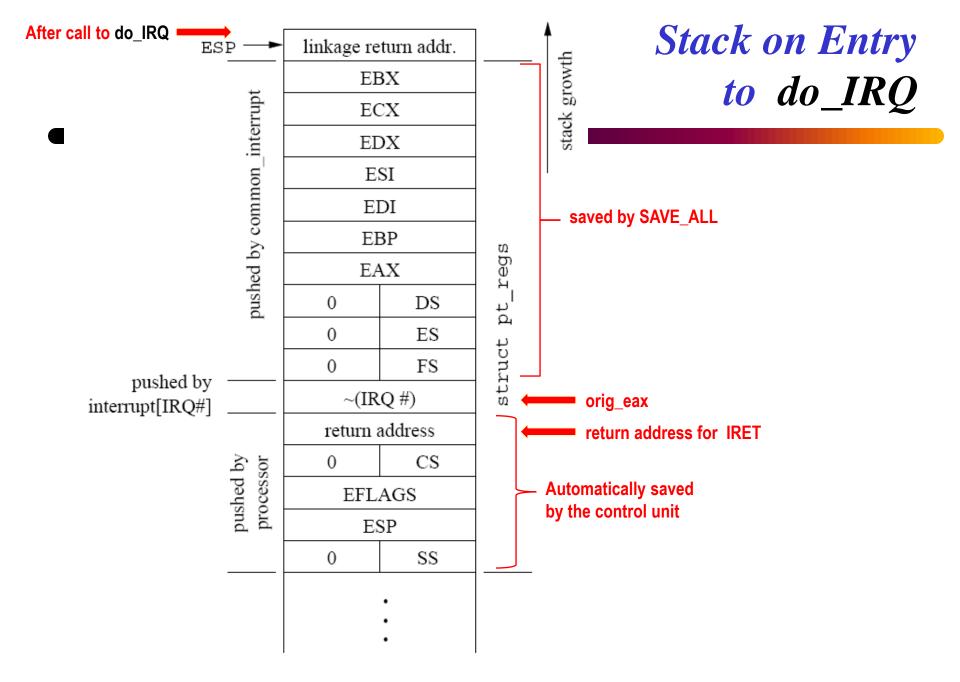
#### MP3 Teams

- Please submit this google form: https://forms.gle/FfeSuDFjo49VNjtd7 before Thursday 3/11 at 11:59PM CT.
- Every team must have exactly 4 members
- Only one person from each group must submit the form

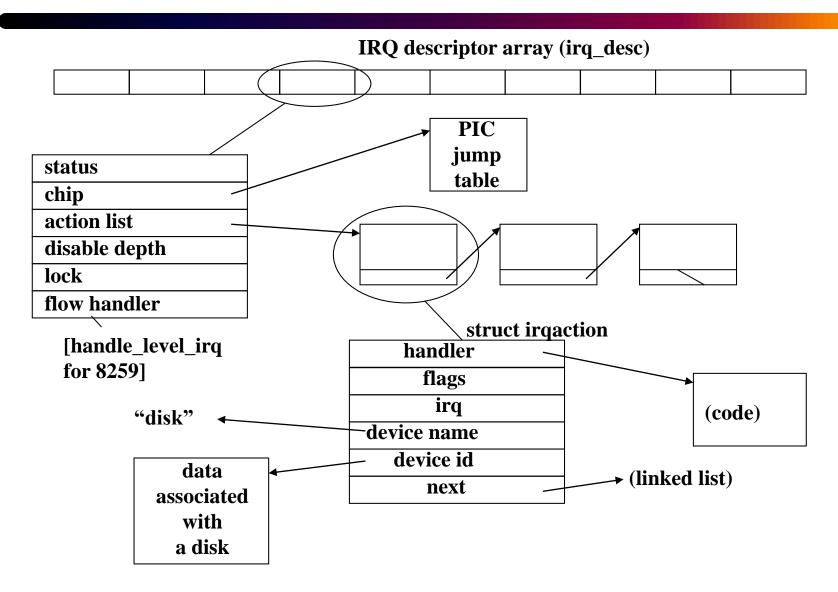
MP3 will be released next week

#### Interrupt Invocation





## Linux' Interrupt Data Structures



# 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 irg (irg, desc); | call interrupt flow handler handle_level_irg 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
    spin unlock (&desc->lock);
```

#### 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
- 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

#### Interrupt Descriptor Table (IDT)

- Associates the interrupt line with the int. handler routine.
- 256 entries (each 8-bytes) or descriptors; each corresponds to an interrupt vector
  - hardware interrupts mapped into vectors 0x20 to 0x2F
- All Linux interrupt handlers are activated by so called: interrupt gates (a descriptor type)
- Whenever an interrupt gate is hit, interrupts are disabled automatically by the processor

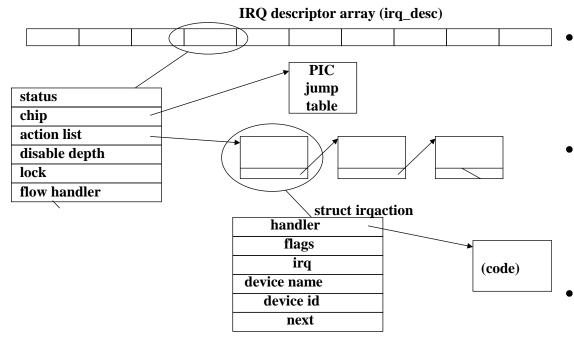
IDT
interrupt[x]

 Before the kernel enables interrupts it must initialize the *idtr* register to point to the IDT table (set by the kernel using lidt instruction)

#### IDT Initialization

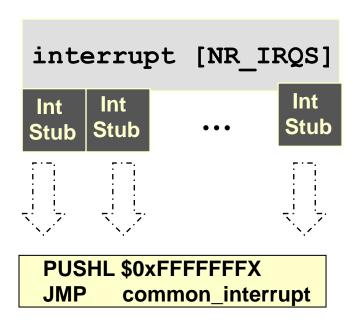
Kernel initialization
 (setup\_idt()) fills
 all the 256 entries of
 IDT with the
 provisional (or null)
 handler

## irq\_descArray



- Every interrupt vector has its own irq\_desc\_t descriptor
- Descriptors are grouped together in irq\_desc array, a data structure supported by Linux
- When a device driver calls the request\_irq() function a new structure to represent the handler is allocated and initialized

#### interrupt[NR IRQS] Array



 Kernel maintains one global array of function pointers
 (interrupt[NR\_IRQS]) in which it stores pointers to interrupt stubs
 (NR\_IRQS is 16 if we use the PIC)

## Initialization of Interrupt Gates

- During initialization init\_IRQ() sets the status field of each IRQ descriptor to IRQ\_DISABLED
- init\_IRQ() updates the IDT by replacing the provisional interrupt gates with new ones

```
for (i = 0; i < NR_IRQS; i++)
    if (i+32 != 128)
        set_intr_gate(i+32, interrupt[i]);</pre>
```

 Interrupt gates are set to the addresses found in the interrupt [NR\_IRQS] array