ECE391 Computer System Engineering Lecture 10

Dr. Zbigniew Kalbarczyk
University of Illinois at Urbana- Champaign

Fall 2018

Lecture Topics

- Linux abstraction of PIC
- General interrupt abstractions
- Linux interrupt system
 - data structures
 - handler installation & removal
 - invocation
 - execution
 - tasklets

ECE391 EXAM 1

- EXAM I October 2; 7:00pm-9:00pm
 - Seating:
 - ECEB 1002: last name starting with A to W
 - ECEB 1013: last name starting with X to Z

Conflict Exam

- Tuesday October 2, 2:00pm
- Location: ECEB 2015
- NO Lecture on Tuesday, October 2
- Review Session
 - Saturday September 29, 4:00pm
 - Location: ECEB 1002

Z. Kalbarczyk

Linux Abstraction of PICs

- Uses a jump table
 - same as vector table (array of function pointers)

- Table is hw_irq_controller structure (or struct irq_chip)
 - each vector # associated
 with a table
 - table used to interact with appropriate PIC (e.g., 8259A, or Advanced PIC)

human-readable name
startup function
shutdown function
enable function
disable function
mask function
mask_ack function
unmask function
(+ several others)

Linux Abstraction of PICs

- hw_irq_controller structure definition
 - IRQs are #'d 0-15 (correspond to vector # 0x20)

```
const char* name;
unsigned int (*startup) (unsigned int irq);
void (*shutdown) (unsigned int irq);
void (*enable)...

void (*disable)...

void (*ack)...

void (*end)...
/* we'll ignore the others... */
```

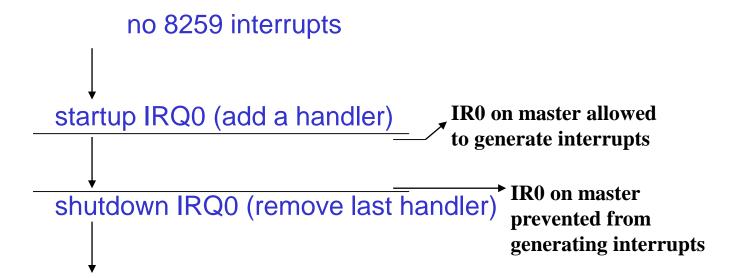
PIC Functions in Jump Table: Explanation

• Initially, all 8259A interrupts are masked out using mask on 8259A

- startup and shutdown functions
 - startup is called when first handler is installed for an interrupt
 - shutdown is called after last handler is removed for an interrupt
 - both functions change the corresponding mask bit in 8259A implementaion

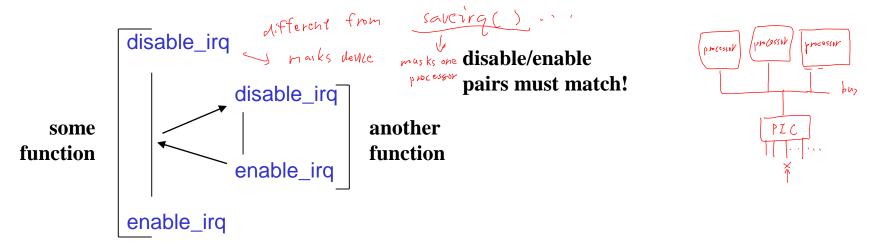
PIC Functions in Jump Table: Explanation (cont.)

Example



PIC Functions in Jump Table (cont.)

- disable/enable functions
 - used to support nestable interrupt masking (disable_irq, enable_irq)



- on 8259
 - first disable_irq calls jump table disable, which masks interrupt on PIC
 - last enable_irq calls jump table enable, which unmasks interrupt on PIC

PIC Functions in Jump Table (cont.)

ack function

acknowledge merrl

- called at start of interrupt handling to ack receipt of the interrupt
- on 8259 (mask and ack), masks interrupt on PIC, then sends EOI to PIC
- end function
 - called at end of interrupt handling
 - on 8259, enables interrupt (unmasks it) on PIC

```
on take an interrupt IS (interrupt to service) set to 1

uck function

EDI

MASK

Grenains!

end function

MASK

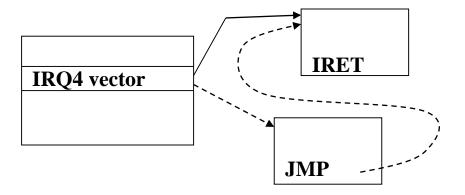
Grenains!
```

General Interrupt Abstractions: Interrupt Chaining

- Hardware view: 1 interrupt → 1 handler
- Problems
 - may have > 15 devices
 - > 1 software routines may want to act in response to device
 - examples:
 - hotkeys for various functions
 - move mouse to lower-right corner to start screen-saver

General Interrupt Abstractions: Interrupt Chaining (cont.)

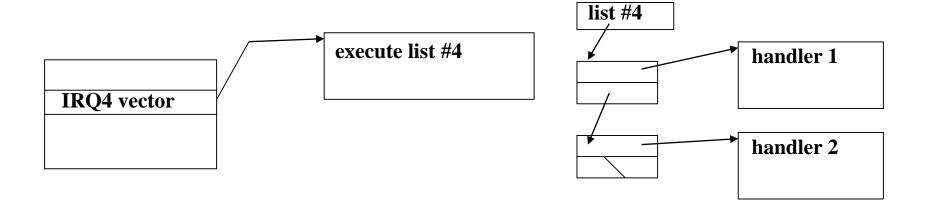
- One approach
 - used by terminate and stay resident (TSR) programs in DOS
 - form linked list (chain) of handlers using JMP instructions
 - not very clean
 - · no way to remove self
 - unless you're first in list
 - to be fair
 - TSR program not designed for removal



General Interrupt Abstractions Interrupt Chaining (cont.)

Solution

- interrupt chaining with linked list data structure
- (not list embedded into code!)



General Interrupt Abstractions: Interrupt Chaining (cont.)

- Drawbacks of chaining
 - for > 1 device
 - must query devices to see if they raised interrupt
 - not always possible
 - for 1 device
 - must avoid stealing data/confusing device
 - example
 - by sending two characters to serial port
 - in response to interrupt declaring port ready for one char.
 - another example
 - reading mouse location twice
 - if device protocol specifies reading once per interrupt

General Interrupt Abstractions: Soft Interrupts (cont.)

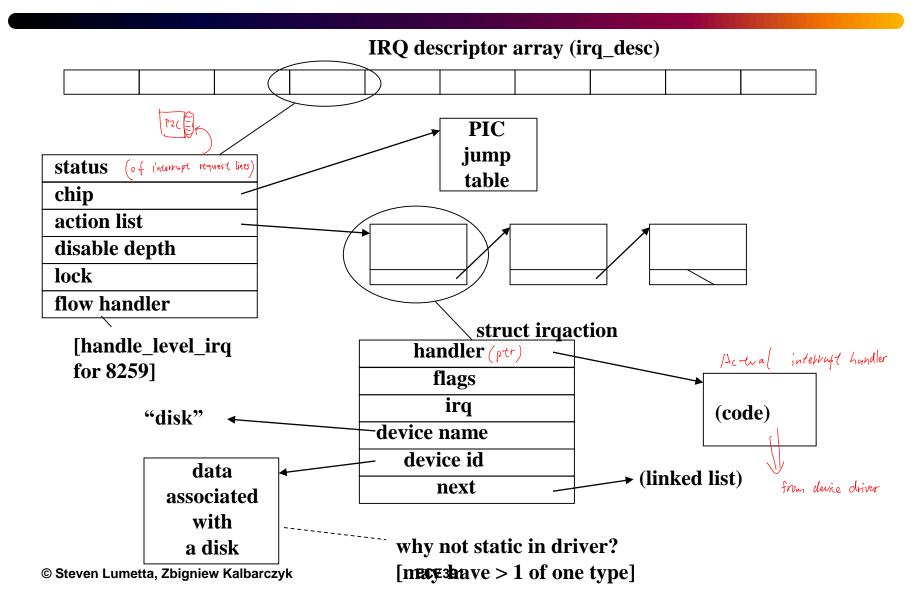
- Recall: why support interrupts?
 - slow device gets timely attention from fast processor
 - processor gets device responses without repeatedly asking for them
- A useful concept in software
 - example: network encryption/decryption
 - packet arrives, given to decrypter
 - when decrypter (software program) is done
 - want to interrupt program
 - to transfer data from packet
 - but has no access to INTR pin

General Interrupt Abstractions: Soft Interrupts (cont.)

Solution

- software-generated (soft) interrupt
- (similarly, but later, signals—user-level soft interrupts)
- runs at priority between program and hard interrupts
- usually generated
 - by hard interrupt handlers
 - to do work not involving device
- Linux version is called tasklets
 - used by code provided to you for MP1
 - discussed later

Linux' Interrupt Data Structures



I inux' request_irq int irq # (0-15 for PIC) request irq (unsigned int irq, (kernel/irq/manage.c) pointer to interrupt handler irq handler t handler, unsigned long irqflags, bit vector const char* devname, human-readable device name void* dev id) pointer to device-specific data struct irgaction* action; int retval; if (irq >= NR IRQS) return -EINVAL; if (!handler) return -EINVAL; action = kmalloc (sizeof (struct irqaction), GFP ATOMIC); - Dynamically allocate if (!action) new action structure for return -ENOMEM; linked list of actions handler action->handler = handler; flags Fill in the new action->flags = irqflags; irq cpus clear (action->mask); action structure device name action->name = devname; device id action->next = NULL; next action->dev id = dev id; select smp affinity (irq); Add the new action structure retval = setup irq (irq, action) to the link list of actions if (retval) kfree (action); Upon failure free the action structure return retval;

Comments on request_irq - Arguments

- **irq irq** + (0-15 for PIC)
- handler pointer to interrupt handler with following arguments
 - irq #
 - dev_id pointer (type void*)
 - handlers should return 1 if handled, 0 if not
- irqflags bit vector
 - IRQF_SHARED interrupt chaining is allowed
 - IRQF_DISABLED execute handlers with IF=0
 - IRQF_SAMPLE_RANDOM use device timing as source of random numbers
- devname human-readable device name (visible in /proc/interrupts)
- dev id pointer to device-specific data (returned to handler when called)

Comments on request_irq (cont.)

- Start by checking values (two checks from the shown code)
- Dynamically allocate new action structure for linked list of actions
- Fill in the new action structure
- Try to add it
 - using setup_irq
 - if call fails, free the action structure

```
int
     setup irq (unsigned int irq, struct irqaction* new)
         struct irq desc* desc = irq desc + irq;
                                                                                         Linux setup_irq
         struct irgaction* old;
         struct irqaction** p;
         unsigned long flags;
         int shared = 0;
                                                                                                    (kernel/irq/manage.c)
         if (irq >= NR IRQS)
             return -EINVAL;
         if (desc->chip == &no irq chip)
             return -ENOSYS;
         if (new->flags & IRQF SAMPLE RANDOM)
             rand initialize irq (irq);
                                                                           Critical section begins
         spin lock irgsave (&desc->lock, flags);
         p = &desc->action;
         if ((old = *p) != NULL) {;
             /* Can't share interrupts unless both agree to and are same type. */
             if (!((old->flags & new->flags) & IRQF SHARED) ||
                 ((old->flags ^ new->flags) & IRQF TRIGGER MASK)) {
                 spin unlock irgrestore (&desc->lock, flags);
                 return -EBUSY;
             /* add new interrupt at end of irq queue */
                                                                            New action is added at the end of the link list
                 p = &old->next;
                 old = *p;
             } while (old);
                                                                                Interrupt
                                                                                              status
             shared = 1;
                                                                                              chip
                                                                                descriptor
                                                                                              action list
         *p = new;
                                                                                              disable depth
         if (!shared) {
                                                                                                                               PIC
             irq chip set_defaults (desc->chip);
                                                                                              lock
             desc->status &= "(IRQ AUTODETECT | IRQ WAITING | IRQ INPROGRESS);
                                                                                              flow handler
                                                                                                                               jump table
             if (!(desc->status & IRQ_NOAUTOEN)) {
                 desc->depth = 0;
                                                                                                                   human-readable name
                 desc->status &= "IRO DISABLED:
                                                                 Startup PIC for this interrupt
                                                                                                                   startup function
                 desc->chip->startup (irq);
                                                                                                                   shutdown function
             else
                                                                                                                   enable function
                 /* Undo nested disables: */
                                                                                                                   disable function
                 desc->depth = 1;
                                                                                                                   mask function
                                                                                                                   mask ack function
                                                                                                                   unmask function
                                                                         Critical section ends
         spin unlock irgrestore (&desc->lock, flags);
                                                                                                                   (+ several others...)
         new->irg = irg:
                                                                  Create new /proc/irg/<irg#> directory
         register irq proc (irq);
         new->dir = NULL:
                                                                  Add handler subdirectory in /proc/irq/<irq#>/<action name>
         register handler proc (irq, new);
© Steven Lumetta, Zbigniew Kalbarczyk
                                                           ECE391
```

Comments on stup_irq

Start with a couple of sanity checks

- If random sampling flag is set
 - initialize as source of random numbers
 - good random numbers are hard to find!
 - devices such as disks can be used because of "random" rotational latency to read data (for example)

Comments on stup_irq (cont.)

- Critical section (most of function)
 - blocks other activity on descriptor
 - uses irqsave/restore to allow handlers to be added from any context

 Note that new action goes at the <u>end</u> of the linked list, not the start

- Interrupt chaining only allowed
 - if <u>all</u> handlers agree to it
 - otherwise only first handler is successfully added

Comments on stup irq (cont.)

- If this handler is the first
 - make sure that PIC jump table has proper default functions
 - clear some status flags
 - clear any previous software disablement (depth)
 - call the PIC startup function

- After critical section
 - create /proc/irq/<irq#> directory if necessary
 - add handler subdirectory in /proc/irq/<irq #>/<action name>

```
void
free irq (unsigned int irq, void* dev id)
                                                                                Linux free_irq
    struct irq desc* desc;
    struct irgaction** p;
                                                                                  (kernel/irq/manage.c)
    unsigned long flags;
    if (irq >= NR IRQS)
        return;
                                                     Find the correct irq descriptor
    desc = irq desc + irq;
    spin lock irqsave (&desc->lock, flags);
                                                           Critical section begins
    p = &desc->action;
    for (;;) {
        struct irgaction* action = *p;
        if (action) {
                                                  Search for action in list
            struct irgaction** pp = p;
            p = &action->next;
                                                    The "continue" causes the rest of the statement
            if (action->dev id != dev id)
                                                    body in the loop to be skipped.
                continue:
            /* Found it - now remove it from the list of entries */
            *pp = action->next;
            if (!desc->action) {
                                                              No more handlers for the IRQ:
                desc->status |= IRQ DISABLED;
                                                               you can invoke PIC shutdown
                desc->chip->shutdown (irq);
            spin unlock irqrestore (&desc->lock, flags);
                                                                      Critical section ends
             /* Remove from /proc/irg/<irg #> */
                                                             Remove handler subdirectory in /proc/irg/<irg#>/<action name>
            unregister handler proc (irq, action);
             /* Make sure it's not being used on another CPU */
            synchronize_irq (irq);
                                Free the action structure
            kfree (action);
            return;
        printk (KERN ERR "Trying to free already-free IRQ %d\n", irq);
                                                                              Critical section ends
        spin unlock irqrestore (&desc->lock, flags);
        return:
     Steven Lumetta, Zbigniew Kalbarczyk
                                                    ECE391
```

Comments on free_irq

- Arguments
 - **irq irq** # (0-15 for PIC)
 - dev_id MUST match pointer value used in request_irq call

Start by checking arguments

- Critical section starts to prevent manipulations of descriptor
 - blocks interrupt from <u>starting</u> to execute
 - (interrupt may <u>already</u> be executing on another processor)

Comments on free_irq (cont.)

- Search for action in list
 - based on dev id pointer
 - hence need for matching pointer
 - passing NULL dev id not allowed with chained interrupts
- Once found (notice that most of the loop executes exactly once)
 - remove the action
 - if this action/handler was the last for this interrupt,
 - turn on software disablement
 - (doing so signals interrupts waiting for descriptor lock to abort once they obtain the lock)
 - call the PIC shutdown function

Comments on free_irq (cont.)

• Remove the /proc/irq/<irq #>/<action name> entry

- On an SMP
 - interrupt may be executing on another processor
 - need to wait for it to finish after releasing descriptor lock

Finally, free the action structure