# Operating System Design and Implementation
## Process Management – Part II

Shiao-Li Tsao

---

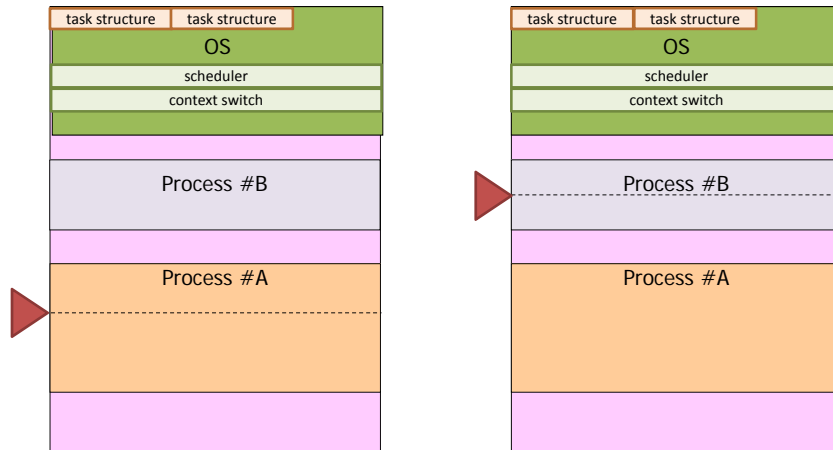| | |
|---|---|
| CPU runs Process A | Execute |
| System call | Schedule |
| Program, process, thread | Context switch |
| Fork | CPU runs Process B |

---

## Process schedule and context switching in Linux

- Scheduling
  - Find the next suitable process to run
- Context switch
  - Store the context of the current process, restore the context of the next process
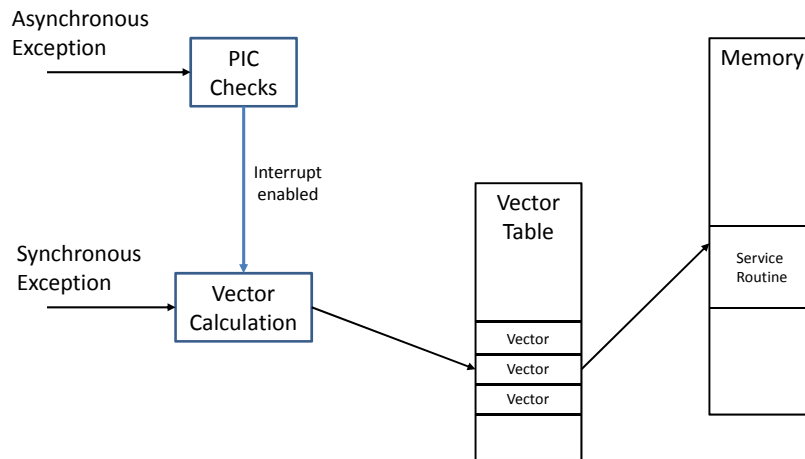
4

# Scheduler in Details

| task structure | task structure |
|---|---|
| OS | |
| scheduler | |
| context switch | |

Process #B

Process #A

---

| task structure | task structure |
|---|---|
| OS | |
| scheduler | |
| context switch | |

Process #B

Process #A

---

# Process schedule and context switching in Linux

- When is the scheduler be invoked
  - Direct invocation vs. Lazy invocation

  - When returning to user-space from a system call
  - When returning to user-space from an interrupt handler
  - When an interrupt handler exits, before returning to kernel-space
  - If a task in the kernel explicitly calls schedule()
  - If a task in the kernel blocks (which results in a call to schedule())

6

# Interrupt Basics

Asynchronous Exception

Synchronous Exception

PIC Checks

Interrupt enabled

Vector Calculation

Vector Table

Vector

Vector

Vector

Memory

Service Routine

*Source: Qing Li "real-time concepts for embedded systems"*

# X86 Interrupts

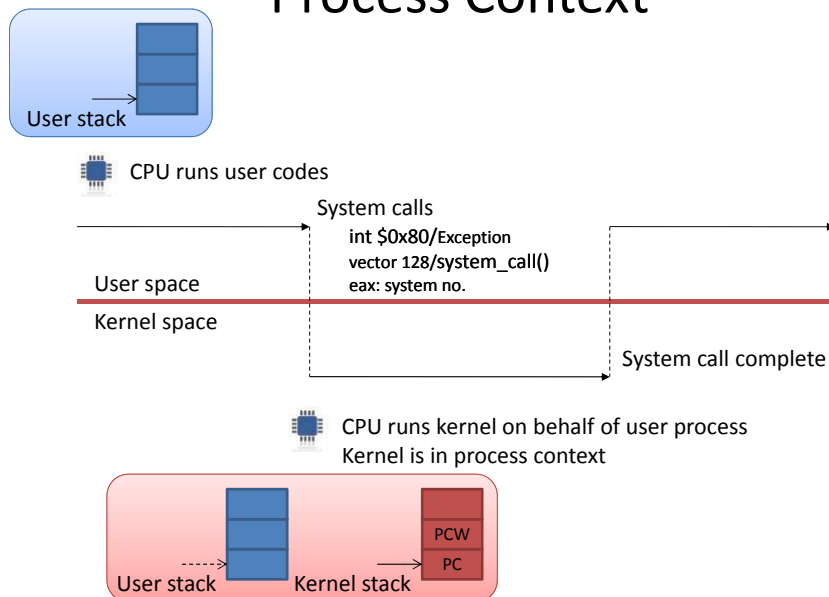| Vector range | Use |
|---|---|
| 0x0 – 0x13 | Non-maskable interrupts/exceptions |
| 0x14 – 0x1F | Intel-reserved |
| 0x20 – 0x7F | External interrupts (IRQs) |
| 0x80 | Programmed exception for system calls |
| 0x81 – 0xEE | External interrupts (IRQs) |

| IRQ | Function |
|---|---|
| 0 | System Timer |
| 1 | Keyboard Controller |
| 2 | 2nd IRQ Controller Cascade |
| 8 | Real-Time Clock |
| 9 | Avail. |
| 10 | Available |
| 11 | Available |
| 12 | Mouse Port / Available |
| 13 | Math Coprocessor |
| 14 | Primary IDE |
| 15 | Secondary IDE |
| 3 | Serial 2 |
| 4 | Serial 1 |
| 5 | Sound card / Parallel 2 |
| 6 | Floppy Disk Controller |
| 7 | Parallel 1 |

*Source: https://en.wikipedia.org/wiki/Interrupt_request_(PC_architecture)*

# Time Interrupt Basics

| OS |
| --- |
| scheduler |
| context switch |
| timer_ISR |
| Vector table |

(4)
(5)
(3)
(5')
(2)

| Process #B |
| --- |

| Process #A |
| --- |

(1)

# Process Context

User stack

CPU runs user codes

System calls
int $0x80/Exception
vector 128/system_call()
eax: system no.

User space

Kernel space

System call complete

CPU runs kernel on behalf of user process
Kernel is in process context

User stack     Kernel stack     PCW     PC

# Interrupt Context

User stack

CPU runs user codes

Interrupt occurs

User space

Kernel space

Interrupt handler returns

CPU runs interrupt handler (or called ISR) in the kernel space
Kernel is in interrupt context

User stack    Interrupt stack
PCW
PC

# Interrupt Context

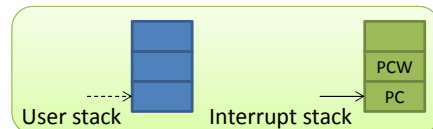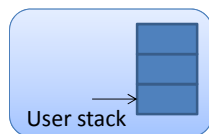User stack

CPU runs user codes

System calls

User space

Kernel space

Interrupt handler returns

CPU runs interrupt handler (or called ISR) in the kernel space
Kernel is in interrupt context

PCW
PC
PCW
PC
Interrupt stack

When returning to user-space from a system call

## If a task in the kernel blocks (which results in a call to schedule())



Usr

AP1   AP2   AP3

Lib1   Lib2

Knl

Syscall

Schr   Klib1   Klib2

Drv1   Drv2

HW

## When returning to user-space from an interrupt handler



Usr

AP1   AP2   AP3

Lib1   Lib2

Knl

Syscall

Schr   Klib1   Klib2

Drv1   Drv2

HW

When returning to user-space from an interrupt handler



When an interrupt handler exits, before returning to kernel-space

## If a task in the kernel explicitly calls schedule()

AP1　　AP2　　AP3

Lib1　　Lib2

Usr

Syscall

Schr　sys_sleep　Klib2

Drv1　　Drv2

Knl

HW

---

# Process context + interrupt context

→ Process #1 context

→ Process #2 context

→ Interrupt #1 context

❶ AP calls read()

❷ read() lib

❸ read() syscall

User space

Kernel space

❹ syscall handler

❺ VFS/buffer cache

❻ FS

❼ Device driver

❽ Set to sleep

❾ Schedule()

❿ Signal process #1

# User preemption

- User preemption occurs when the kernel is in a safe state and about to return to user-space

User process #2

User process #1

User space

Kernel space

Interrupt handler

Call schedule()
Check need_resched
About to return to user space

# User preemption

- User preemption occurs when the kernel is in a safe state and about to return to user-space

User process #2

User process #1

User space

Kernel space

System call

Call schedule()
Check need_resched
About to return to user space

# Kernel preemption

- Linux kernel is possible to preempt a task at any point, so long as the kernel does not hold a lock

User process #2

User process #1

User space

Kernel space

System call

Interrupt handler

Call schedule()
Check need_resched
About to return to user space

# Preemptive Kernel

- Non-preemptive kernel supports user preemption
- Preemptive kernel supports kernel/user preemption

- Kernel can be interrupted ≠ kernel is preemptive
    - Non-preemptive kernel, interrupt returns to interrupted process
    - Preemptive kernel, interrupt returns to any schedulable process

# Preemptive Kernel

- 2.4 is a non-preemptive kernel
- 2.6 is a preemptive kernel
- 2.6 could disable CONFIG_PREEMPT

# Preemptive Kernel

# Preemptive Kernel

- How difficult to implement a preemptive kernel?



User process #2

User process #1

User space

Kernel space

*Kernel code must be reentrant*

System call #1

Interrupt handler

System call #1

Call schedule()
Check need_resched
About to return to user space

---

# Single Core vs. Multi-core



CPU State
Interrupt Logic
Execution Units | Cache

Single processor/single core

CPU State
Interrupt Logic
Execution Units | Cache

CPU State
Interrupt Logic
Execution Units | Cache

Multiple processor/single core

CPU State
Interrupt Logic

CPU State
Interrupt Logic

Execution Units | Cache

Single processor/single core/hyper threading

CPU State
Interrupt Logic
Execution Units | Cache

CPU State
Interrupt Logic
Execution Units | Cache

Single processor/multi core/separated cache

CPU State
Interrupt Logic
Execution Units

CPU State
Interrupt Logic
Execution Units

Cache

Single processor/multi core/shared cache

CPU State
Interrupt Logic

CPU State
Interrupt Logic

Execution Units | Cache

CPU State
Interrupt Logic

CPU State
Interrupt Logic

Execution Units | Cache

Single processor/multi core/hyper threading

# Single Core vs. Multi-core

| CPU State | CPU State | CPU State | CPU State |
|---|---|---|---|
| Interrupt Logic | Interrupt Logic | Interrupt Logic | Interrupt Logic |

| Execution Units | Cache | Execution Units | Cache |

Linux Kernel

Main Memory

shared resources for all cores (critical section)

shared resources per core (critical section)

---

# Schedule Algorithms

- Think about yourself (your homework schedule)
  - Given that you have a lot of homework to do, each with a deadline
  - Profs. continue assigning new homework

  - What is the next homework to do? (next task to schedule)
  - Why should we stop a homework? (time to schedule)
  - How long can we concentrate on a homework? (scheduling period)
  - How long do we spend to determine the next homework? (scheduling algorithm overhead)
  - How much effort do we spend to switch homework? (context switch overhead)
  - What is the importance of a homework? (priority of a job)
  - How long does a homework need? (job length)

# How Linux Scheduler Works

Priority = static priority (nice) + dynamic priority (heuristic)

Time slice = Func(static priority)

PD #2

From parent's

**Nice (-20 to +19): normal tasks OR**
Real-time priority (0 to 99) : real-time tasks

Process list | PD #1 | PD #2 | PD #3 | PD #4 | PD #5 | PD #6

Run queue

---

# Timeslice

- Timeslice function

| Type of Task | Nice Value | Timeslice Duration |
|---|---|---|
| Initially created | parent's | half of parent's |
| Minimum Priority | +19 | 5ms (MIN_TIMESLICE) |
| Default Priority | 0 | 100ms (DEF_TIMESLICE) |
| Maximum Priority | -20 | 800ms (MAX_TIMESLICE) |

lower priority or less interactive ←        → higher priority or more interactive

Minimum 5ms    Default 100ms    Maximum 800ms

$$\text{base time quantum (in milliseconds)} = \begin{cases} (140 - static\ priority) \times 20 & \text{if } static\ priority < 120 \\ (140 - static\ priority) \times 5 & \text{if } static\ priority \geq 120 \end{cases} \quad (1)$$

## Process schedule and context switching in Linux

- Priority-based scheduler
- Dynamic priority-based scheduling
  - Dynamic priority
    - Normal process
      - nice value: -20 to +19 (larger nice values imply you are being nice to others)
  - Static priority
    - Real-time process
      - 0 to 99
  - Total priority: 140

33

## Linux O(1) cheduler



schedule()

sched_find_first_set()

bit 7 (priority 7)

bit 0 (priority 0)

list of all runnable tasks, by priority

140-bit priority array     bit 139 (priority 139)

Run the first process in the list

# Linux O(1) cheduler

schedule()

sched_find_first_set()

bit 0 (priority 0)

bit 7 (priority 7)

**active**

list of all runnable
tasks, by priority

**expired**

**Time_slice = 0**

140-bit priority array

bit 139 (priority 139)

**Timeslice = Func(static priority)**

Run the first process in the list


# Calculating Priority

- static_prio = nice
- Prio = nice – bonus + 5

  dynamic priority = max (100, min ( static priority - bonus + 5, 139))

- Heuristic
  - sleep_avg: (0 to MAX_SLEEP_AVG(10ms))
  - sleep_avg+=sleep (becomes runnable)
  - Sleep_avg-=run (every time tick when task runs)

# System calls related to scheduling

| System call | Description |
|---|---|
| nice( ) | Change the static priority of a conventional process |
| getpriority( ) | Get the maximum static priority of a group of conventional processes |
| setpriority( ) | Set the static priority of a group of conventional processes |
| sched_getscheduler( ) | Get the scheduling policy of a process |
| sched_setscheduler( ) | Set the scheduling policy and the real-time priority of a process |
| sched_getparam( ) | Get the real-time priority of a process |
| sched_setparam( ) | Set the real-time priority of a process |
| sched_yield( ) | Relinquish the processor voluntarily without blocking |
| sched_get_ priority_min( ) | Get the minimum real-time priority value for a policy |
| sched_get_ priority_max( ) | Get the maximum real-time priority value for a policy |
| sched_rr_get_interval( ) | Get the time quantum value for the Round Robin policy |
| sched_setaffinity( ) | Set the CPU affinity mask of a process |
| sched_getaffinity( ) | Get the CPU affinity mask of a process |

# How Linux Scheduler Works



SCHED_FIFO or SCHED_RR

Priority = static priority (nice)

Time slice = N/A (FIFO)
Time slice = config

PD #2

Nice (-20 to +19): normal tasks OR
**Real-time priority (0 to 99) : real-time tasks**

Process list | PD #1 | PD #2 | PD #3 | PD #4 | PD #5 | PD #6

Run queue

The following flow boxes appear in the slide:

- CPU runs Process A
- Execute
- System call
- Schedule
- Program, process, thread
- Context switch
- Fork
- CPU runs Process B

---

# Process schedule and context switching in Linux

- Context switch
  - Hardware context switch
    - Task State Segment Descriptor (Old Linux)
  - Step by step context switch
    - Better control and optimize
- Context switch
  - switch_mm()
    - Switch virtual memory mapping
  - switch_to()
    - Switch processor state
- Process switching occurs only in kernel mode
- The contents of all registers used by a process in User Mode have already been saved

40

# Scheduler in Details

Task state segment

| TSS (A) | TSS (B) |
|---|---|

OS

scheduler

context switch

Process #B

Process #A

LDR
Registers
...

IP
SP
EAX
EBX

CPU

*TSS structure source: http://www.embedded.com/design/prototyping-and-development/4025054/Managing-Tasks-on-x86-Processors*

---

# Scheduler in Details

Task state segment

| TSS (A) | TSS (B) |
|---|---|

OS

scheduler

context switch

Process #B

Process #A

LDR
Registers
...

IP
SP
EAX
EBX

CPU

*TSS structure source: http://www.embedded.com/design/prototyping-and-development/4025054/Managing-Tasks-on-x86-Processors*

# Scheduler in Details

Task state segment

| TSS (A) | TSS (B) |
|---------|---------|

OS

scheduler

context switch

Process #B

Process #A

LDR
Registers
...

IP
SP
EAX
EBX

CPU

*TSS structure source: http://www.embedded.com/design/prototyping-and-development/4025054/Managing-Tasks-on-x86-Processors*



# Scheduler in Details

Task state segment

| TSS (A) | TSS (B) |
|---------|---------|

OS

scheduler

context switch

Process #B

Process #A

LDR
Registers
...

IP
SP
EAX
EBX

CPU

*TSS structure source: http://www.embedded.com/design/prototyping-and-development/4025054/Managing-Tasks-on-x86-Processors*

# Scheduler in Details

| TSS (A) | TSS (B) | |
|---|---|---|

**Task state segment**

**OS**

scheduler

context switch

Process #B

Process #A

**CPU**
- IP
- SP
- EAX
- EBX

LDR
Registers
...

*TSS structure source: http://www.embedded.com/design/prototyping-and-development/4025054/Managing-Tasks-on-x86-Processors*

---

# How x86 helps in Context Switch

Kernel code seg.

GDT

Kernel data seg.

LDT #1

TSS #1

LDT #2

Task 2 data seg.

Task 2 code seg.

Kernel code seg.

GDT

Kernel data seg.

LDT #1

TSS #1

LDT #2

Task 1 code seg.

Task 1 data seg.

**GDT**
- LDT k
- TSS k
- ...
- LDT 1
- TSS 1
- sys
- Kernel Data
- Kernel Text

**LDT**
- NULL
- Task k .text
- Task k .data
- ...

TSS
structure

# Process related terms



CPU → registers

Storage
program

process

Memory

process

Process
Related info

Code segment

Data segment

Stack segment

*Is that good enough ?*
*If not, why ?*

Physical memory might be discontinuous

47

---

# Process related terms (Cont.)



CPU → registers   registers   registers

Storage
program

Process related info

process

Thread dedicate pages

Thread dedicate pages

Thread dedicate pages

thread related info

thread related info

thread related info

Code segment

Data segment

Stack segment

Code segment

Data segment

Stack segment

Code segment

Data segment

Stack segment

48

# Process related terms (Cont.)

- Depending on OS designs

| | | |
|---|---|---|
| **User** | **User** | **User** |
| **Kernel** | **Kernel** | **Kernel** |
| Processes only | User thread | kernel thread |

49

# Process related terms (Cont.)

- Linux lightweight process

Processes share many resources and pages

**User**

**Kernel**

50