


# Writing the operating system from scratch



# Questions 1: how to boot?

What we Know:

BIOS loads first sector from boot device into memory 0x7c00, then jump to it.

Solutions:

Write [a boot program](#) into a boot device.

Choose [a floppy disk](#) as a boot device.

Boot program size limited by 512 bytes, so make [a filesystem](#) on floppy disk, write [a loader program](#) into it.

Boot program loads this loader program into memory, then jump to it.

Loader program is designed to implement more complex functions

# Environment



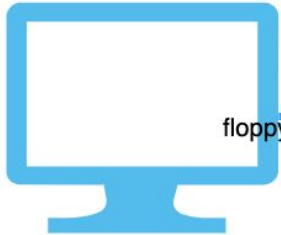
boot.bin



Loader.bin



kernel.bin



bochs

floppya: 1\_44=image/my\_os.img, status=inserted



my\_os.img

## Install OS

```
dd if=build/boot.bin of=image/my_os.img bs=512 count=1 conv=notrunc
```

```
sudo mount image/my_os.img /mnt/floppy
```

```
sudo cp build/loader.bin build/kernel.bin /mnt/floppy
```

```
sleep 1 sudo
```

```
umount /mnt/floppy
```

## Make a floppy disk image

```
bximage
```

```
mkfs.msdos my_os.img
```

# Boot

Read a sector from **Root Directory**

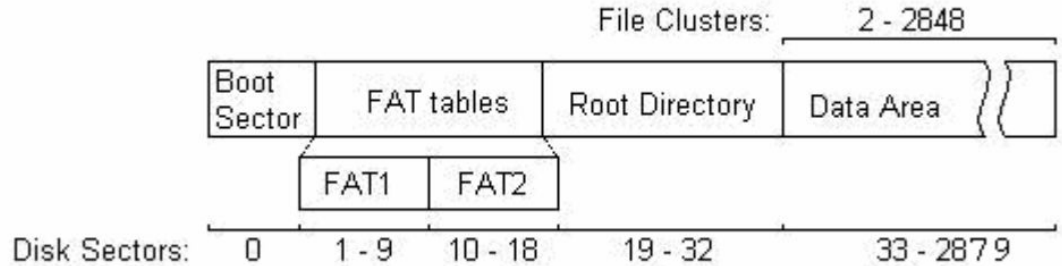
For each directory entry(32 Bytes)

Find name "loader.bin"

Get **First Logical Cluster** No.

Read sector from **Data Area** into memory 0x90000

JMP 0x9000:0x100



each entry(12 bits) in **FAT tables** pointer to a cluster/sector on **Data Area**

Specifically, the FAT entry values signify the following:

Value	Meaning
0x00	Unused
0xFF0-0xFF6	Reserved cluster
0xFF7	Bad cluster
0xFF8-0xFFF	Last cluster in a file
(anything else)	Number of the next cluster in the file

# Questions 2: how to manage memory?

What we Know:

Real mode is characterized by a 20-bit segmented memory address space (giving exactly 1 MiB of addressable memory) and unlimited direct software access to all addressable memory, I/O addresses and peripheral hardware.

Solutions:

Let's talk about Protected Mode

# Segmentation

Segmentation provides a mechanism for dividing the processor's addressable memory space (called the **linear address** space) into smaller protected address spaces called segments.

**Segments** can be used to hold the code, data, and stack for a program.

To locate a byte in a particular segment, a **logical address** (also called a far pointer) must be provided. A logical address consists of a **segment selector** and an **offset**

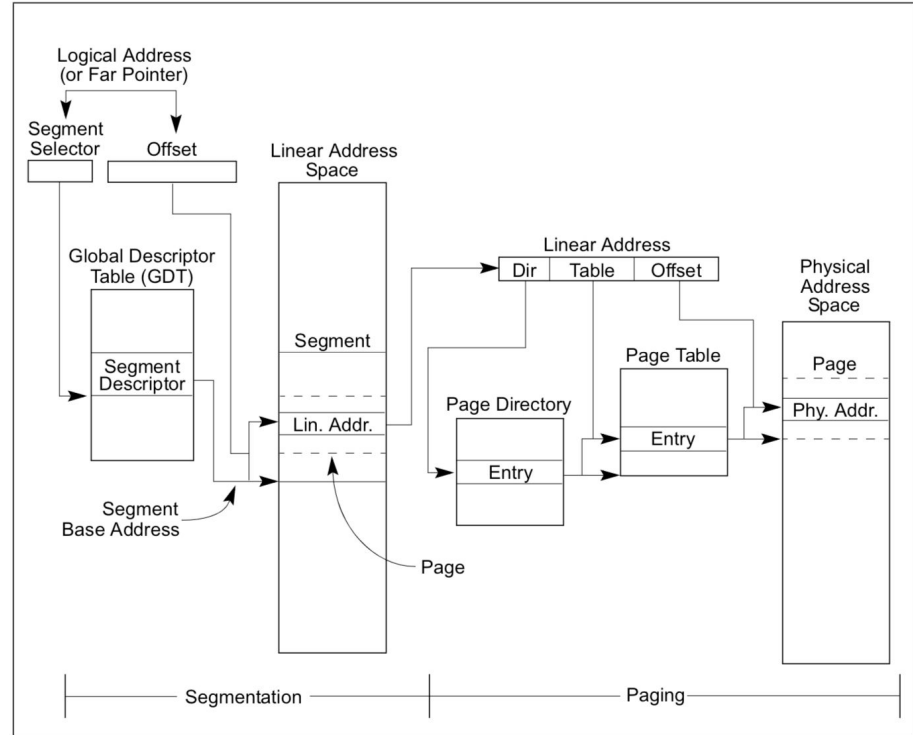


Figure 3-1. Segmentation and Paging

# Segment Selector

## Segment Descriptor

A [segment selector](#) is a 16-bit identifier for a segment (see Figure 3-6). It does not point directly to the segment, but instead points to the segment descriptor that defines the segment.

A [segment descriptor](#) is a data structure in a GDT or LDT that provides the processor with the size and location of a segment, as well as access control and status information.

Segment descriptors are typically created by compilers, linkers, loaders, or the operating system or executive, but not application programs.

```

41  %macro Descriptor 3
42  →   dw %2 & 0FFFFh →   →
43  →   dw %1 & 0FFFFh →   →
44  →   db (%1 >> 16) & 0FFh →   →
45  →   dw ((%2 >> 8) & 0F00h) | (%3 & 0F0FFh) →
46  →   db (%1 >> 24) & 0FFh →   →
47  %endmacro

```

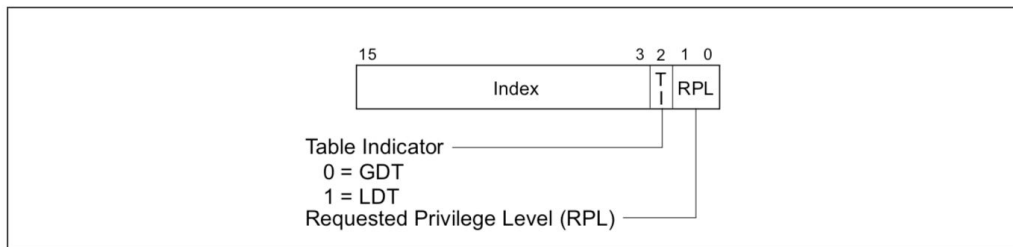
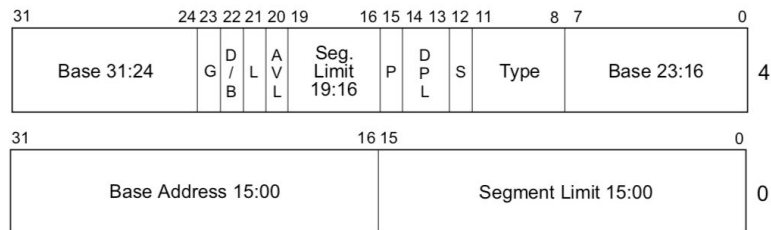


Figure 3-6. Segment Selector



- L — 64-bit code segment (IA-32e mode only)
- AVL — Available for use by system software
- BASE — Segment base address
- D/B — Default operation size (0 = 16-bit segment; 1 = 32-bit segment)
- DPL — Descriptor privilege level
- G — Granularity
- LIMIT — Segment Limit
- P — Segment present
- S — Descriptor type (0 = system; 1 = code or data)
- TYPE — Segment type

Figure 3-8. Segment Descriptor

# GDT & LDT

The GDT is not a segment itself; instead, it is a data structure in linear address space. The base linear address and limit of the GDT must be loaded into the GDTR register.

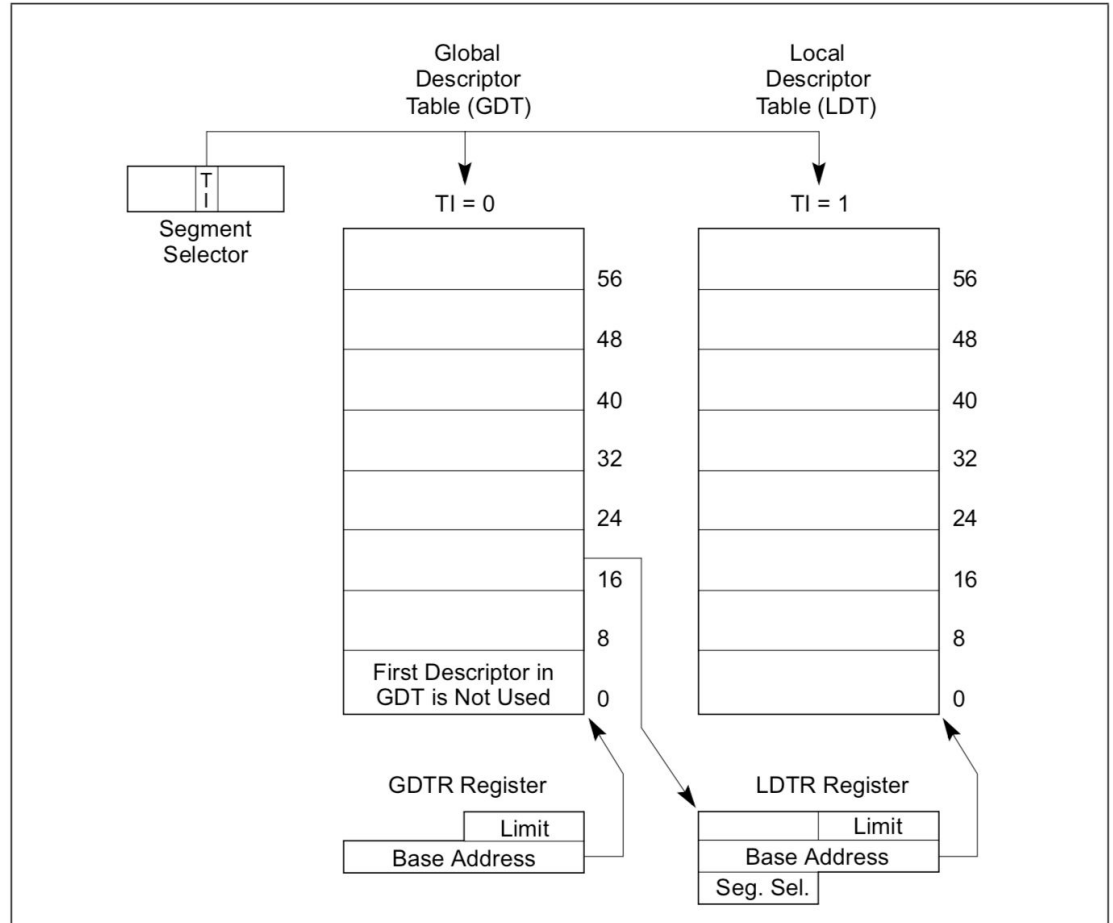


Figure 3-10. Global and Local Descriptor Tables



# Codes

Prepare GDT

Set cr0 PE flag

Initialize segment register

```
9      , GDT
10     LABEL_GDT: Descriptor 0, 0, 0
11     LABEL_DESC_FLAT_C: Descriptor 0, 0fffffh, DA_CR|DA_32|DA_LIMIT_4K; 0-4G
12     LABEL_DESC_FLAT_RW: Descriptor 0, 0fffffh, DA_DRW|DA_32|DA_LIMIT_4K; 0-4G
13     LABEL_DESC_VIDEO: Descriptor 0b8000h, 0ffffh, DA_DRW|DA_DPL3;
14
15     GdtLen equ $ - LABEL_GDT
16     GdtPtr dw GdtLen - 1; boundary
17     ... dd BaseOfLoader*10h + LABEL_GDT
18     ; Selector
19     SelectorFlatC equ LABEL_DESC_FLAT_C - LABEL_GDT
20     SelectorFlatRW equ LABEL_DESC_FLAT_RW - LABEL_GDT
21     SelectorVideo equ LABEL_DESC_VIDEO - LABEL_GDT + SA_RPL3
```

```
153     ... lgdt [GdtPtr]
154     ... cli
155
156     ... in al, 92h
157     ... or al, 00000010b
158     ... out 92h, al
159
160     ... mov eax, cr0
161     ... or eax, 1
162     ... mov cr0, eax
163
164     ... jmp dword SelectorFlatC:(BaseOfLoaderPhyAddr+LABEL_PM_START)
```

```
292 LABEL_PM_START:
293     ... mov ax, SelectorVideo
294     ... mov gs, ax
295
296     ... mov ax, SelectorFlatRW
297     ... mov ds, ax
298     ... mov es, ax
299     ... mov fs, ax
300     ... mov ss, ax
```

# Paging

Paging provides a mechanism for implementing a conventional demand-paged, virtual-memory system where sections of a program's execution environment are mapped into physical memory as needed.

Paging is enabled if **CR0.PG = 1**. Paging can be enabled only if protection is enabled (CR0.PE = 1).

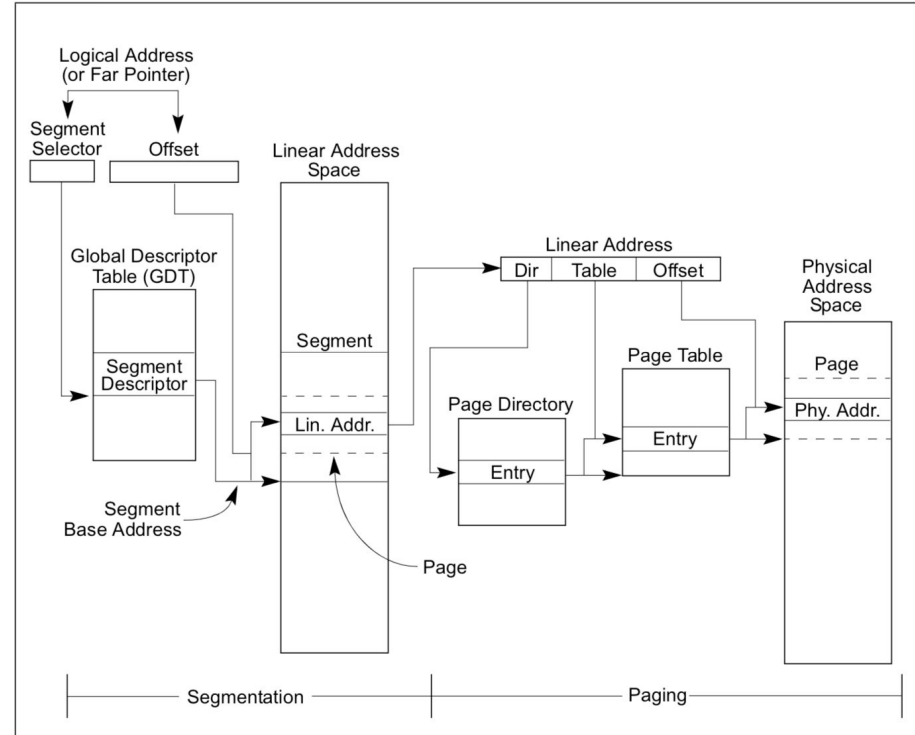


Figure 3-1. Segmentation and Paging

```

366 ; Function SetupPaging
367 SetupPaging:
368     xor edx, edx
369     mov eax, [dwMemSize]
370     mov ebx, 400000h
371     div ebx
372     mov ecx, eax
373     test edx, edx
374     jz .no_remainder
375     inc ecx
376 .no_remainder:
377     push ecx
378     ;
379     ;
380     ;
381     ; init page dir table
382     mov ax, SelectorFlatRW
383     mov es, ax
384     mov edi, PageDirBase
385     xor eax, eax
386     mov eax, PageTblBase | PG_P-- | PG_USU- | PG_RWW
387 .1:
388     stosd
389     add eax, 4096
390     loop .1
391     ;
392     ; init page table
393     pop eax
394     mov ebx, 1024
395     mul ebx
396     mov ecx, eax
397     mov edi, PageTblBase
398     xor eax, eax
399     mov eax, PG_P-- | PG_USU- | PG_RWW
400 .2:
401     stosd
402     add eax, 4096
403     ; page manage 4K space
404     loop .2
405     ;
406     mov eax, PageDirBase
407     mov cr3, eax
408     mov eax, cr0
409     or eax, 80000000h
410     mov cr0, eax
411     jmp short .3
412 .3:
413     nop
414     ret
415

```

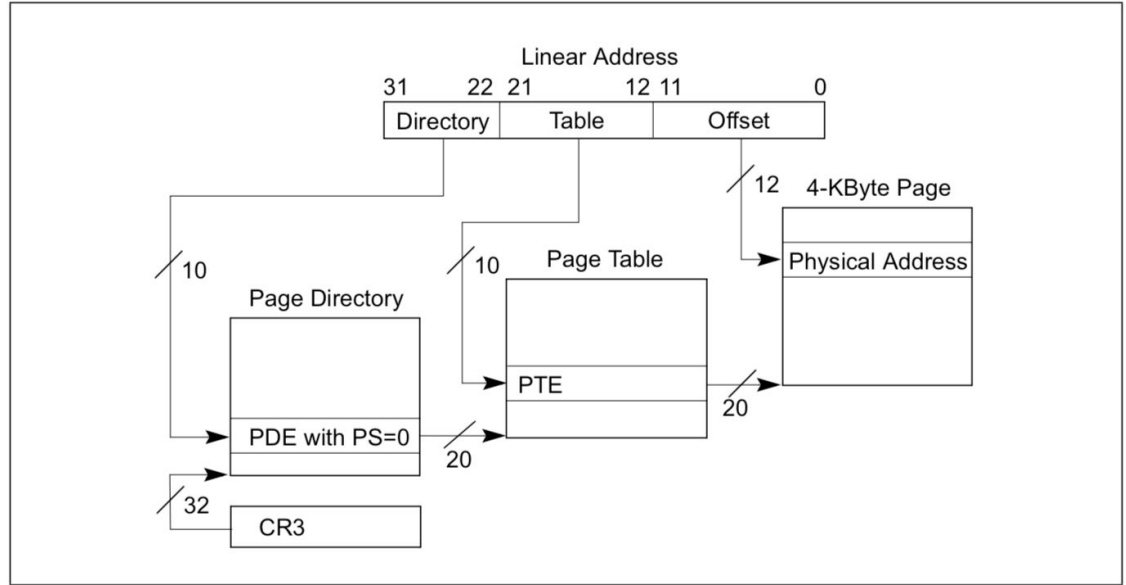


Figure 4-2. Linear-Address Translation to a 4-KByte Page using 32-Bit Paging

# Questions 3: how to design User mode and Kernel mode?

What we Know:

OS like unix usually forbid user accessing hardware directly.

The processor uses privilege levels to prevent a program or task operating at a lesser privilege level from accessing a segment with a greater privilege.

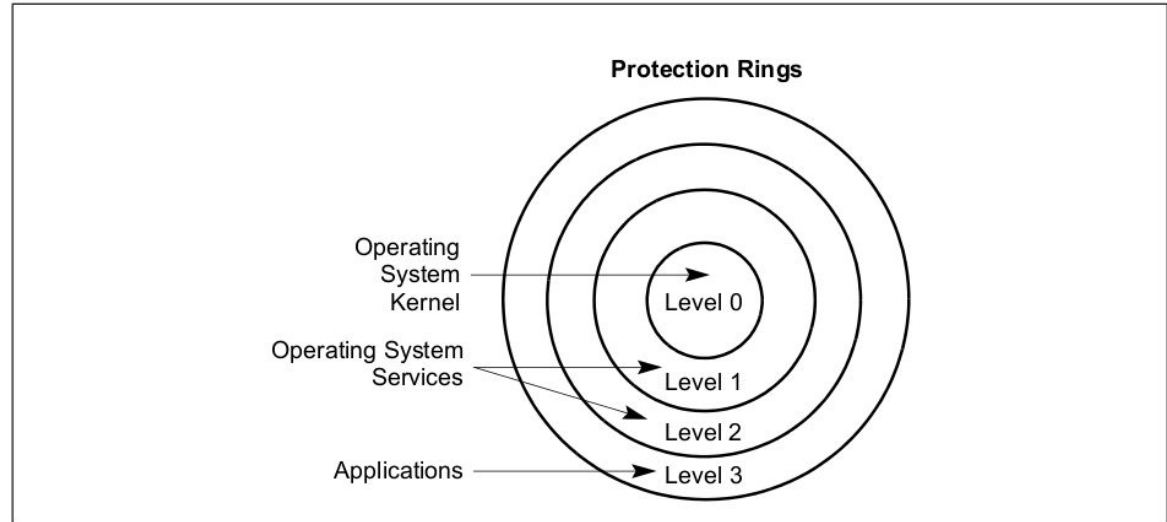
Solutions:

User mode run privilege level 1, implement Process

Kernel mode run privilege level 0, implement Interrupt handle, such as clock to schedule Process

# Privilege Level

The processor's segment-protection mechanism recognizes 4 **privilege levels**, numbered from 0 to 3. The greater numbers mean lesser privileges.



**Figure 5-3. Protection Rings**

# Accessing Data Segments

The **CPL** is the privilege level of the currently executing program or task. It is stored in bits 0 and 1 of the CS and SS segment registers

The **DPL** is the privilege level of a segment or gate. It is stored in the DPL field of the segment or gate descriptor for the segment or gate.

The **RPL** is an override privilege level that is assigned to segment selectors. It is stored in bits 0 and 1 of the segment selector.

$$\text{Max}(\text{CPL}, \text{RPL}) \leq \text{DPL}$$

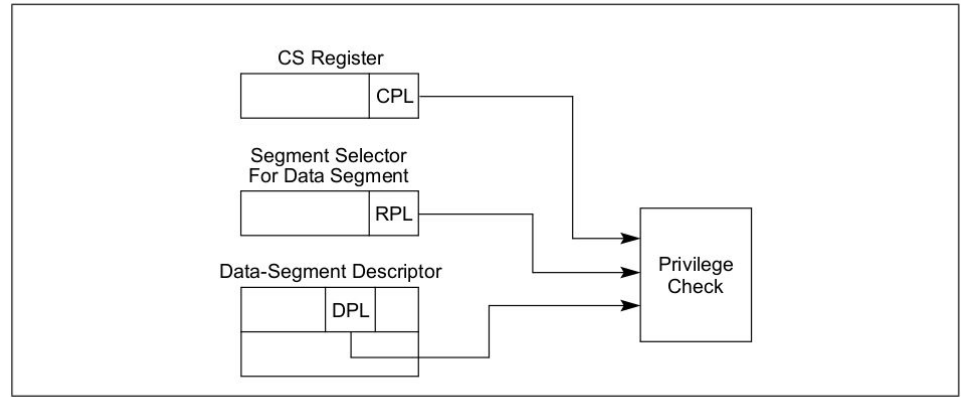


Figure 5-4. Privilege Check for Data Access

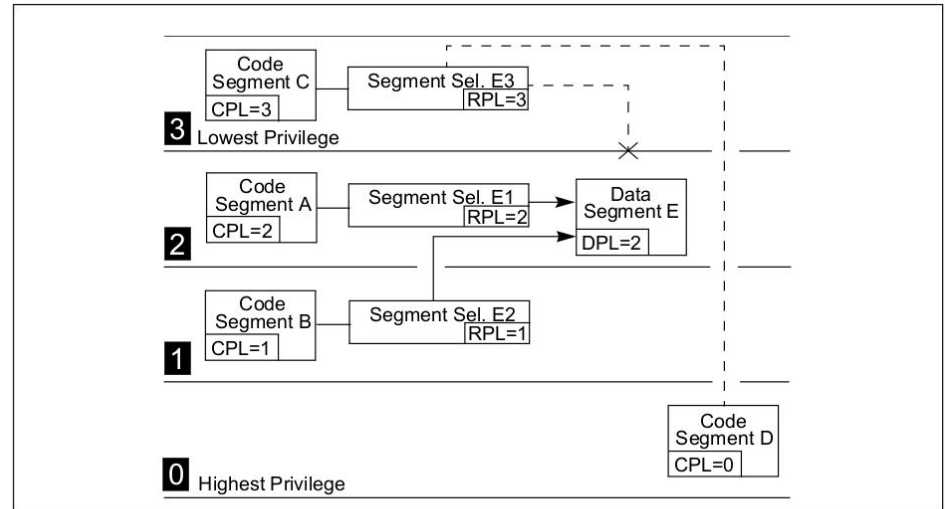


Figure 5-5. Examples of Accessing Data Segments From Various Privilege Levels

# Accessing Code Segments

Most code segments are nonconforming  
For these segments, program control can be transferred only to code segments **at the same level of privilege**, unless the transfer is carried out through a call **gate**, as described in the following sections.

**Nonconforming:**  $CPL = DPL$  and  $RPL \leq DPL$

**Conforming:**  $CPL \geq DPL$ ,  $CPL$  not change

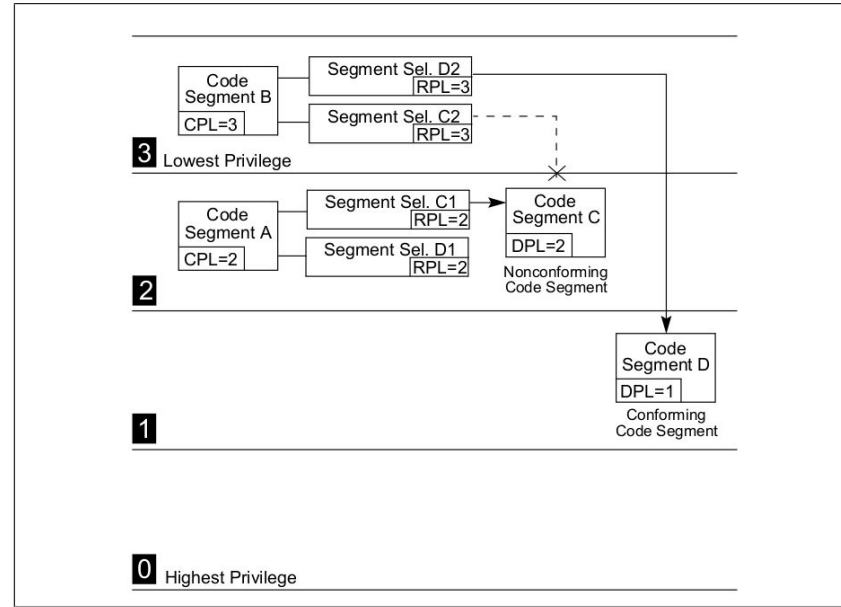


Figure 5-7. Examples of Accessing Conforming and Nonconforming Code Segments From Various Privilege Levels

# Gate Descriptors

To provide controlled access to code segments with **different privilege levels**, the processor provides special set of descriptors called gate descriptors.

There are four kinds of gate descriptors:

- Call gates
- Trap gates
- **Interrupt gates**
- Task gates

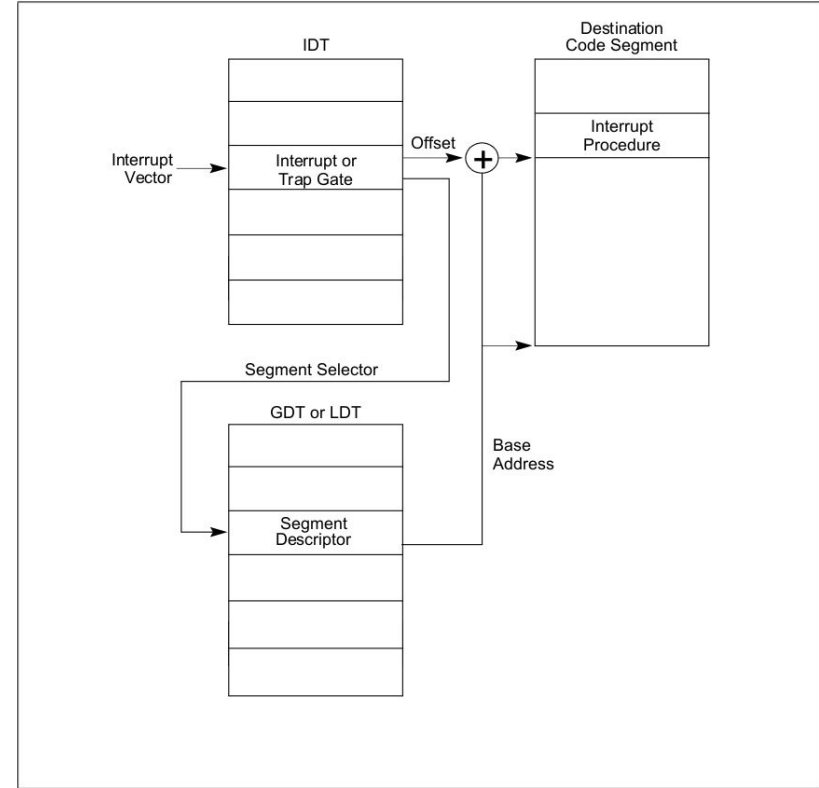
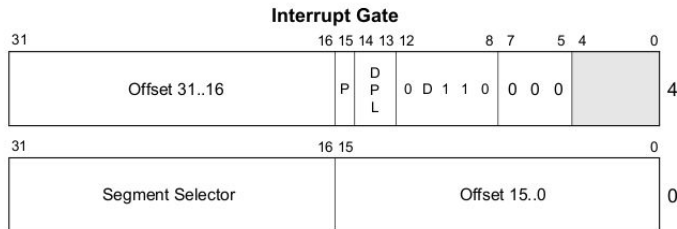


Figure 6-3. Interrupt Procedure Call



# Interrupt Init & Handle

```
216 void init_mask_interrupt()  
217 {  
218     /* master */  
219     init_idt_desc([INT_VECTOR_IRQ0 + 0, DA_386IGate,  
220                 mask_int_func0, PRIVILEGE_KRNL]);  
221  
222     init_idt_desc(INT_VECTOR_IRQ0 + 1, DA_386IGate,  
223                 mask_int_func1, PRIVILEGE_KRNL);
```

```
167 ALIGN 16  
168 mask_int_func0: .....; Interrupt r  
169     mask_int_func_master ..... 0  
170     .....  
171 ALIGN 16  
172 mask_int_func1: .....; Interrupt r  
173     mask_int_func_master ..... 1  
174     .....
```

```
141 %macro mask_int_func_master ..... 1  
142     call save; save reg  
143  
144     in al, INT_M_CTLMASK  
145     or al, (1<<%1)  
146     out INT_M_CTLMASK, al; disable this irq  
147     mov al, EOI  
148     out INT_M_CTL, al  
149  
150     sti; enable other irq  
151  
152     push ..... %1  
153     call ..... [g_irq_table + 4 * %1]  
154     pop ecx  
155  
156     cli; disable all irq  
157  
158     in al, INT_M_CTLMASK  
159     and al, ~(1<<%1)  
160     out INT_M_CTLMASK, al; enable this irq  
161     ret; jmp to restart or restart_reenter  
162 %endmacro  
163 ; -----
```

# Stack Switching

The operating system is responsible for creating **stacks** and stack-segment descriptors for all the privilege levels to be used and for loading initial pointers for these stacks into the **TSS**

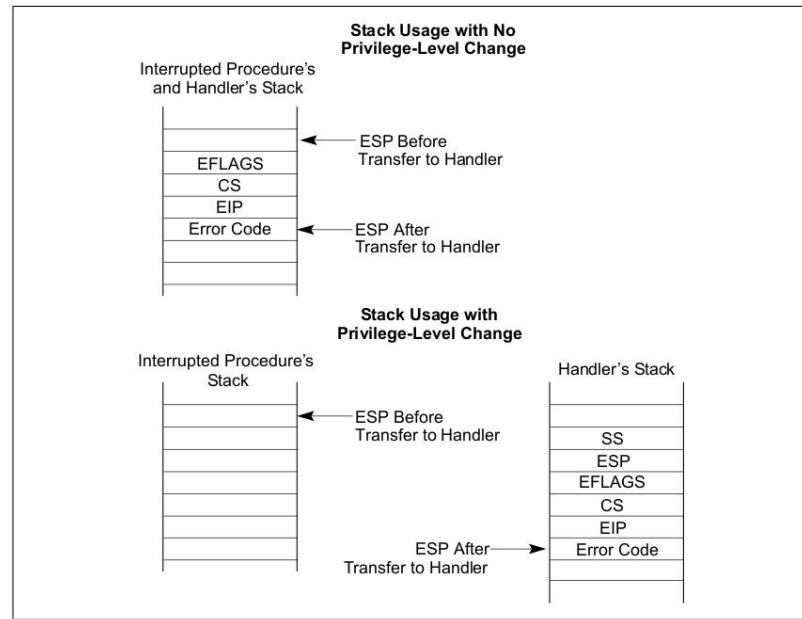


Figure 6-4. Stack Usage on Transfers to Interrupt and Exception-Handling Routines

32 (TSS)		28
Reserved	SS2	24
ESP2		20
Reserved	SS1	16
ESP1		12
Reserved	SS0	8
ESP0		4
Reserved	Previous Task Link	0

Reserved bits. Set to 0.

Figure 7-2. 32-Bit Task-State Segment (TSS)

# Save

Step by step to draw 3 stacks!!!

```
109 ; interrupt occur, save the reg;
110 ; if this irq is the only one, which means
111 ; esp point to process stackframe(not process's stack)
112 ; we need switch to kernel stack
113 ; otherwise, this irq may interrupt other irq,
114 ; stack here is just the kernel stack
115 save:
116     pushad
117     push ds
118     push es
119     push fs
120     push gs
121     mov dx, ss
122     mov ds, dx
123     mov es, dx
124
125     mov esi, esp
126     ....
127     ; Judge reenter or not
128     inc dword [g_k_reenter]
129     cmp dword [g_k_reenter], 0
130     jne .reenter
131     ....
132     ; switch to kernel stack
133     mov esp, StackTop
134     push restart
135     jmp [esi+RETADR-P_STACKBASE]
136 .reenter:
137     push restart_reenter
138     jmp [esi+RETADR-P_STACKBASE]
```

# User Process Restart

Move stackframe address to esp

Move top of stackframe to tss.sp0

Pop ...

Iretq back to process instruction

```
90 ; return from irq,  
91 ; if irq is reenter one, recover the previous irq  
92 ; otherwise, recover the process  
93 restart:  
94     mov esp, [g_proc_ready]  
95     lldt [esp + P_LDT_SEL]  
96     lea eax, [esp + P_STACKTOP]  
97     mov dword [g_tss + TSS3_S_SP0], eax  
98 restart_reenter:  
99     dec dword [g_k_reenter]  
100     pop gs  
101     pop fs  
102     pop es  
103     pop ds  
104     popad  
105     add esp, 4  
106  
107     iretd
```

# Questions 4: how to design Multiprocess?

What we Know:

- Interrupt handle save registers to a stackframe from tss.sp0

- Every process should prepare this stackframe before running

Solutions:

- Define a struct contains stackframe for process

# User Process Data Structure

Prepare stack at Level 0

```
46 extern proc_t* g_proc_ready;
47 extern proc_t g_proc_table[PROC_MAX];
48
49 extern char g_task_stack[STACK_SIZE*PROC_MAX];
50
```

```
8 typedef struct stackframe_s
9 {
10     uint32_t gs;
11     uint32_t fs;
12     uint32_t es;
13     uint32_t ds;
14     uint32_t edi;
15     uint32_t esi;
16     uint32_t ebp;
17     uint32_t kernel_esp;
18     uint32_t ebx;
19     uint32_t edx;
20     uint32_t ecx;
21     uint32_t eax;
22     uint32_t retaddr;
23     uint32_t eip;
24     uint32_t cs;
25     uint32_t eflags;
26     uint32_t esp;
27     uint32_t ss;
28 }stackframe_t;
29
30 typedef struct proc_s
31 {
32     stackframe_t regs;
33     uint16_t ldt_sel;
34     descriptor_t ldts[LDT_SIZE];
35     uint32_t pid;
36     char proc_name[PROC_NAME_MAX];
37 }proc_t;
```

```
2 P_STACKBASE equ 0
3 GSREG equ P_STACKBASE
4 FSREG equ GSREG + 4
5 ESREG equ FSREG + 4
6 DSREG equ ESREG + 4
7 EDIREG equ DSREG + 4
8 ESIREG equ EDIREG + 4
9 EBPREG equ ESIREG + 4
10 KERNELESPPREG equ EBPREG + 4
11 EBXREG equ KERNELESPPREG + 4
12 EDXREG equ EBXREG + 4
13 ECXREG equ EDXREG + 4
14 EAXREG equ ECXREG + 4
15 RETADR equ EAXREG + 4
16 EIPREG equ RETADR + 4
17 CSREG equ EIPREG + 4
18 EFLAGSREG equ CSREG + 4
19 ESPREG equ EFLAGSREG + 4
20 SSREG equ ESPREG + 4
21 P_STACKTOP equ SSREG + 4
22 P_LDT_SEL equ P_STACKTOP
23 P_LDT equ P_LDT_SEL + 4
24
25 TSS3_S_SP0 equ 4
```

# User Process Initialize

Each process has a stack and a entry function.

```
19  proc_t* p_proc = g_proc_table;
20  ....
21  int idx = 0;
22  for (idx = 0; idx < PROC_MAX; ++idx)
23  {
24      p_proc = g_proc_table + idx;
25      p_proc->ldt_sel = SELECTOR_LDT_FIRST + (idx<<3);
26      memcpy(&p_proc->ldts[0], &g_gdt[SELECTOR_KERNEL_CS>>3], sizeof(descriptor_t));
27      p_proc->ldts[0].attr1 = DA_C | PRIVILEGE_TASK << 5; /* change the DPL */
28      memcpy(&p_proc->ldts[1], &g_gdt[SELECTOR_KERNEL_DS>>3], sizeof(descriptor_t));
29      p_proc->ldts[1].attr1 = DA_DRW | PRIVILEGE_TASK << 5; /* change the DPL */
30
31      p_proc->regs.cs = (0 & SA_RPL_MASK & SA_TI_MASK) | SA_TIL | RPL_TASK; /* this RPL is
32      import when iretd */
33      p_proc->regs.ds = (8 & SA_RPL_MASK & SA_TI_MASK) | SA_TIL | RPL_TASK;
34      p_proc->regs.es = (8 & SA_RPL_MASK & SA_TI_MASK) | SA_TIL | RPL_TASK;
35      p_proc->regs.fs = (8 & SA_RPL_MASK & SA_TI_MASK) | SA_TIL | RPL_TASK;
36      p_proc->regs.ss = (8 & SA_RPL_MASK & SA_TI_MASK) | SA_TIL | RPL_TASK;
37      p_proc->regs.gs = (SELECTOR_KERNEL_GS & SA_RPL_MASK) | RPL_TASK;
38      p_proc->regs.eip = (uint32_t)g_func_table[idx];
39      p_proc->regs.esp = (uint32_t)g_task_stack + STACK_SIZE * idx;
40      p_proc->regs.eflags = 0x1202; /* IF=1, IOPL=1, bit 2 is always 1. */
41  }
42
43  g_proc_ready = g_proc_table;
```



# Simple & Stupid Schedule

```
61 void kernel_schedule()  
62 {  
63     ++g_proc_ready;  
64     if (g_proc_ready >= g_proc_table + PROC_MAX)  
65         g_proc_ready = g_proc_table;  
66 }
```



# Reference

<https://software.intel.com/en-us/articles/intel-sdm>