

# **Traps, interrupts, and drivers**

---

Youjip Won



# Contents

---

- System calls, exceptions, and interrupts
- X86 protection
- Code: The first system call
- Code: Assembly trap handlers
- Code: C trap handler
- Code: System calls

# System calls, exceptions, and interrupts

---

- System call
  - When a program asks for an operating system service.
- Exception
  - When a program performs an illegal action.
- Interrupt
  - When a device generates a signal to indicate that it needs attention from the operating system.
  - Kernel handles all interrupts because in most cases, only the kernel has the required privilege and state.

# The first system call: initcode.S

```
exec ("init", ["init", NULL]);
```

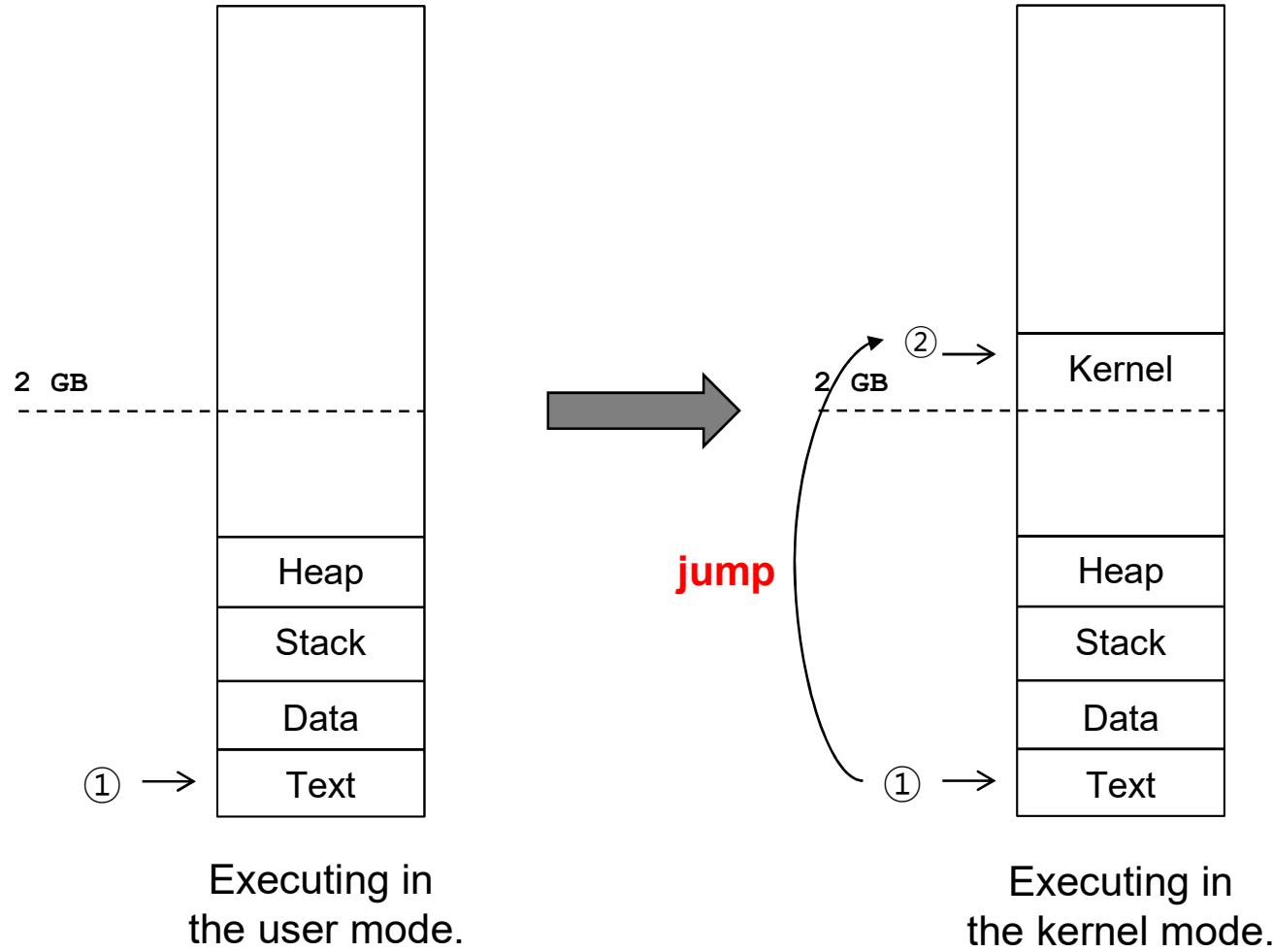
```
1 # Initial process execs /init.  
...  
9 .globl start  
10 start:  
11     pushl $argv  
12     pushl $init  
13     pushl $0 // where caller pc would be  
14     movl $SYS_exec, %eax  
15     int $T_SYSCALL  
...  
...
```

```
...  
23 # char init[] = "/init\0";  
24 init:  
25     .string "/init\0"  
26  
27 # char *argv[] = { init, 0 };  
28 .p2align 2  
29 argv:  
30     .long init  
31     .long 0  
32
```

Switching from Kernel mode to user mode



# Transfer the control from the user to the kernel.



# When does it go into kernel?

---

- ➊ System call
  - ➌ Call 'int' instruction.
  - ➍ `read()`, `fork()`
- ➋ Exception
  - ➌ Generated by the software

```
1 int i,j;  
2 i = 1; j = 0;  
3 i = i/j; !
```

- ➌ Interrupt
  - ➌ Hardware interrupt signal, e.g. ntimer Interrupt.
  - ➌ Only the kernel has the privilege to handle it!

# How does it go into kernel?

---

`int n`

- Assembly instruction
  - It changes the mode from the user to the kernel if necessary.
  - It executes interrupt handler.
- ➊ Interrupt Handler
- A function that handles the interrupt.
- ➋ Find the address of interrupt handler and check the privilege level.
- ➌ Change the stack from the user stack to the kernel stack.
- ➍ Push some value of the registers.
- ➎ Load the address of interrupt handler to `%eip` register.

# int n

---

- ① Find the address of interrupt handler and check the privilege level.
- ② Change the stack from the user stack to the kernel stack.
- ③ Push some value of the registers.
- ④ Load the address of interrupt handler to %eip register.

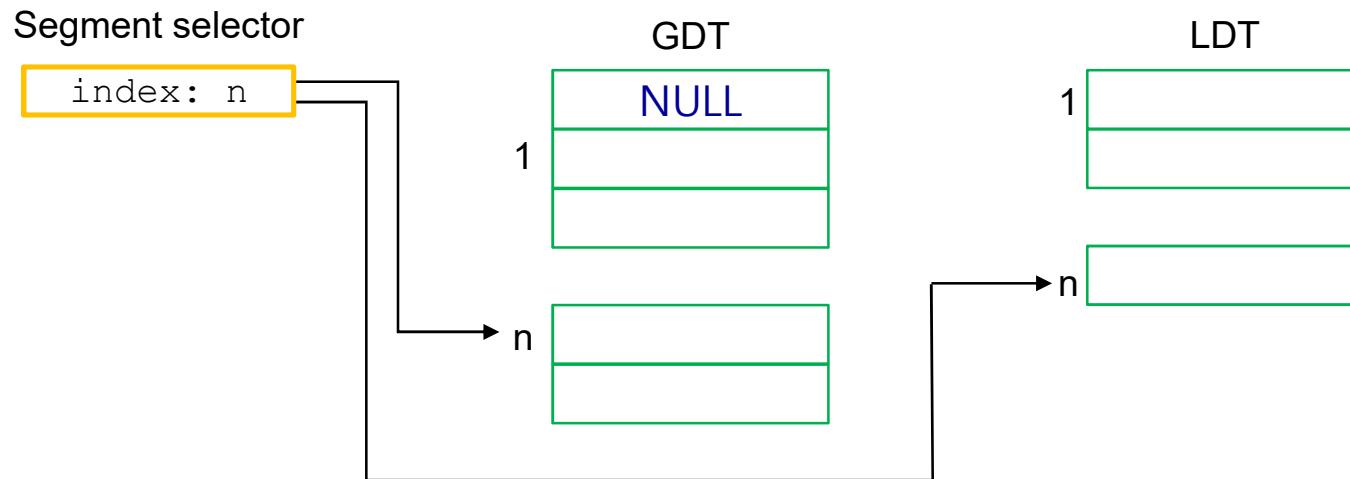
# Privilege Level (Mode)

- Execution Privilege
- Four levels in x86.
  - It is called “ring”
  - 0 (kernel, the height), 1, 2, 3 (user, the lowest)
- **CPL:** current privilege level
  - RPL field (2 bit) of %CS register.



# Recap: Segment Descriptor Table

- Kernel maintains the array of segment descriptors.
  - Per CPU: Global Descriptor Table
  - Per Process: Local Descriptor Table
- Segment Selector (%cs, %es, %ds, %ss, %fs, %gs)
  - It points an entry in a descriptor table.



# Segment Selector



index: index of the descriptor.

TI: table indicator, 0 = GDT, 1 = LDT

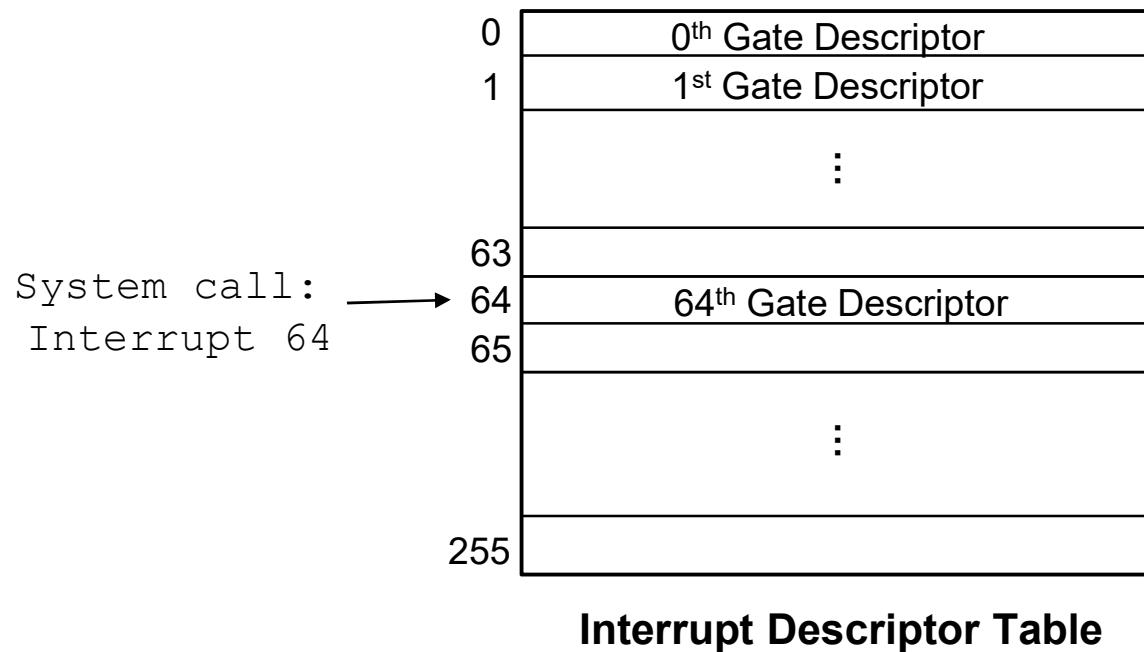
RPL: Request Privilege Level

**Now we can find the current privilege level.**

**What next? Where is the privilege level of interrupt handler?**

# Interrupt Descriptor Table

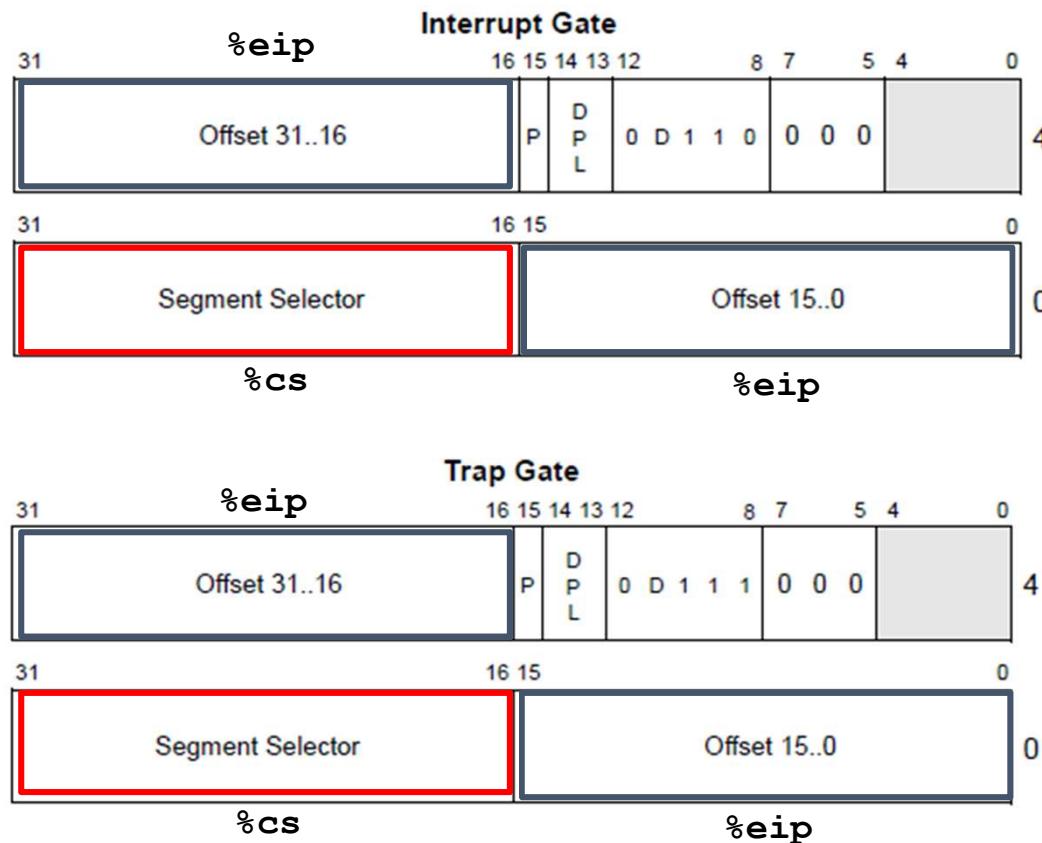
- Kernel maintains the array of **gate descriptors**.
  - Each entry contains the address of interrupt handler.
  - When interrupt N is occurred, the interrupt handler at index N is executed.



# Gate Descriptor

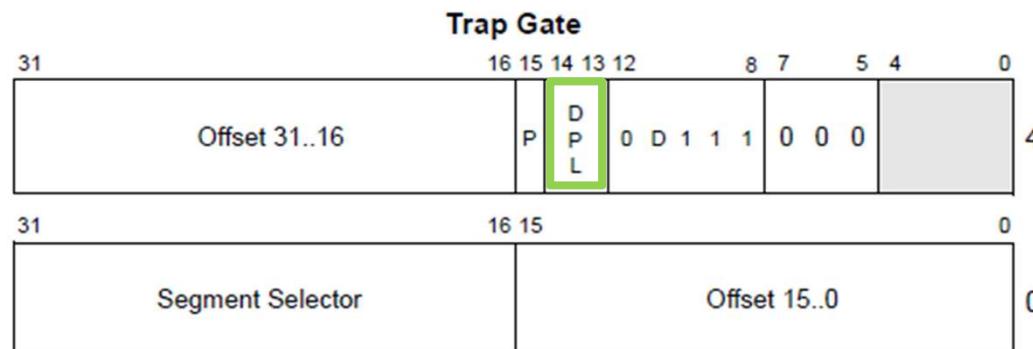
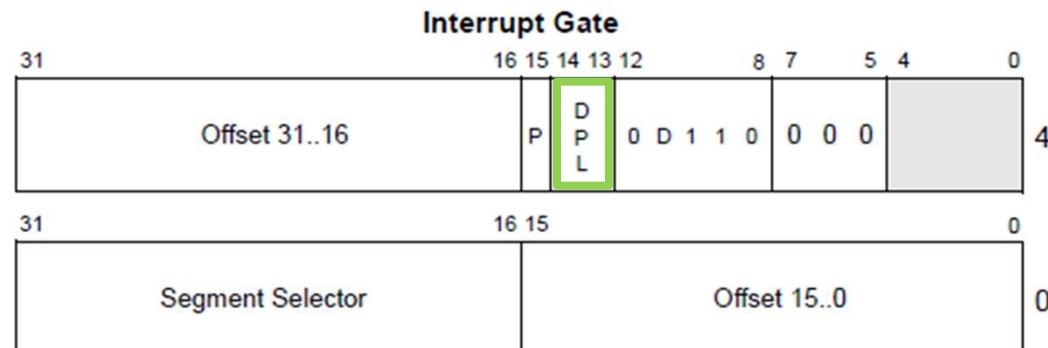
- Special data structure in x86.
- Offset 31..16 & 15..0: Value of %eip register
- Segment Selector: Value of %cs register

Address of interrupt handler



# Gate Descriptor

- Special data structure in x86.
- DPL : Privilege level of interrupt handler.



## struct gate descriptor: mmu.h

```
147 // Gate descriptors for interrupts and traps
148 struct gatedesc {
149     uint off_15_0 : 16;      // low 16 bits of offset in segment
150     uint cs : 16;           // code segment selector
151     uint args : 5;          // # args, 0 for interrupt/trap gates
152     uint rsv1 : 3;          // reserved(should be zero I guess)
153     uint type : 4;          // type(STS_{IG32,TG32})
154     uint s : 1;              // must be 0 (system)
155     uint dpl : 2;            // descriptor(meaning new) privilege level
156     uint p : 1;              // Present
157     uint off_31_16 : 16;     // high bits of offset in segment
158 };
```

# Initializing interrupt descriptor table: main.c

```
17 int
18 main(void)
19 {
    ...
30     tvinit();          // gate descriptors
    ...
35     kinit2(P2V(4*1024*1024), P2V(PHYSTOP));
36     userinit();        // first user process
37     mpmain();          // finish this processor's setup
38 }
```

- ➊ The interrupt descriptor table is initialized when the system is booting up.
  - The main function for kernel calls the function `tvinit()`.
  - The `tvinit()` initializes the interrupt descriptor table.

# Initializing interrupt descriptor table: trap.c

- Initializing 256 entries
- All entries: disable interrupt, can be called only inside the kernel.
- System call: interruptible, can be called from the user (`DPL_USER == 3`).

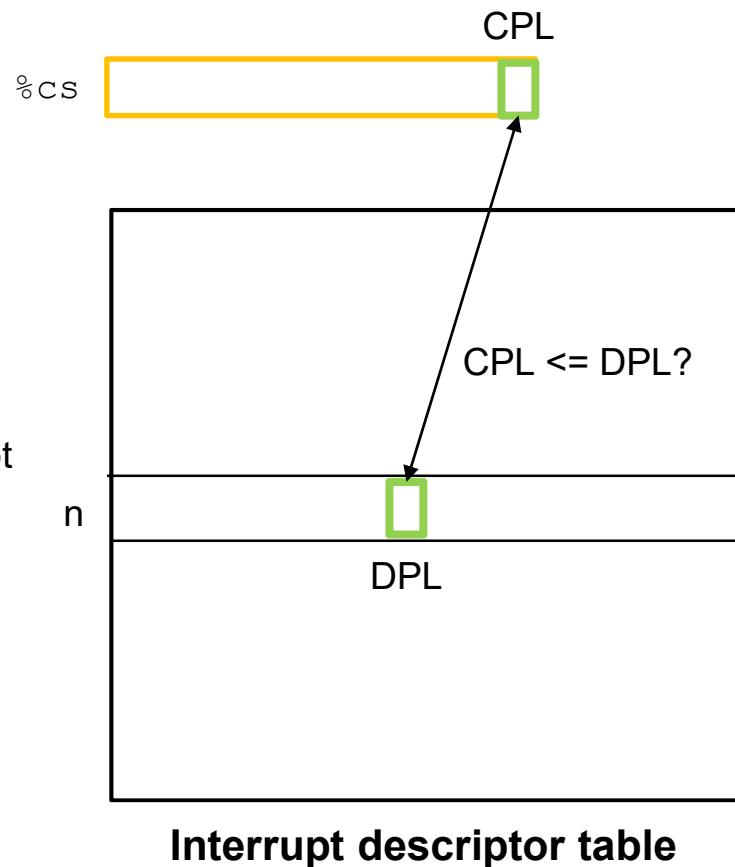
```
12 struct gatedesc idt[256];
...
17 void
18 tvinit(void)
19 {
20     int i;
21
22     for(i = 0; i < 256; i++)
23         SETGATE(idt[i], 0, SEG_KCODE<<3, vectors[i], 0);
24     SETGATE(idt[T_SYSCALL], 1, SEG_KCODE<<3, vectors[T_SYSCALL], DPL_USER);
25
26     initlock(&tickslock, "time");
27 }
```

`SETGATE(Gate descriptor,  
 istrap?, %cs value,  
 %eip value, DPL);`

  
**why?**

# Checking the privilege level

- **CPL:** RPL field of %cs register.
- **DPL:** DPL field of gate descriptor.
- $CPL \leq DPL$ :
  - If CPL is lower (higher), Do next.
  - The execution privilege of current context is higher than the execution privilege of interrupt handler.
- $CPL > DPL$ :
  - Otherwise, Kernel panic!!!
  - Raise Interrupt 13 (T\_GPFLT).



# Checking the privilege level (Cont.)

- The CPL of user mode context is 3.
- Only 64<sup>th</sup> Gate descriptor's DPL is set to be 3 (== DPL\_USER). Others are 0.
- User can call only 'int 64'.
- Hardware interrupt is regardless of whether the execution is in kernel or user context.

```
12 struct gatedesc idt[256];
```

...

```
17 void  
18 tvinit(void)  
19 {  
20     int i;  
21  
22     for(i = 0; i < 256; i++)  
23         SETGATE(idt[i], 0, SEG_KCODE<<3, vectors[i], 0);  
24     SETGATE(idt[T_SYSCALL], 1, SEG_KCODE<<3, vectors[T_SYSCALL],  
25             DPL_USER);  
26     initlock(&tickslock, "time");  
27 }
```

```
SETGATE(Gate descriptor,  
        istrap?, %cs value,  
        %eip value, DPL);
```

DPL\_USER

why?

# int n

---

- ① Find the address of interrupt handler and check the privilege level.
- ② Change the stack from the user stack to the kernel stack.
- ③ Push some value of the registers.
- ④ Load the address of interrupt handler to %eip register.

# Changing the stack

---

- CPL <= DPL:
  - The process can start this point only when the above condition is satisfied.
- CPL == DPL:
  - The privilege level is not updated. The stack does not have to be changed.
- CPL < DPL:
  - The privilege level is updated, stack should be changed from the user stack to the kernel stack.

**How does the processor find the location of the kernel stack?**

# Task State Segment

- Data structure in x86 CPU that holds information on a task.
- Per-CPU data structure

```
2 struct cpu {  
3     uchar apicid;           // Local APIC ID  
4     struct context *scheduler; // swtch() here to enter scheduler  
5     struct taskstate ts;      // Used by x86 to find stack for interrupt  
6     struct segdesc gdt[NSEGS]; // x86 global descriptor table  
7     volatile uint started;    // Has the CPU started?  
8     int ncli;                // Depth of pushcli nesting.  
9     int intena;              // Were interrupts enabled before pushcli?  
10    struct proc *proc;        // The process running on this cpu or null  
11 };
```

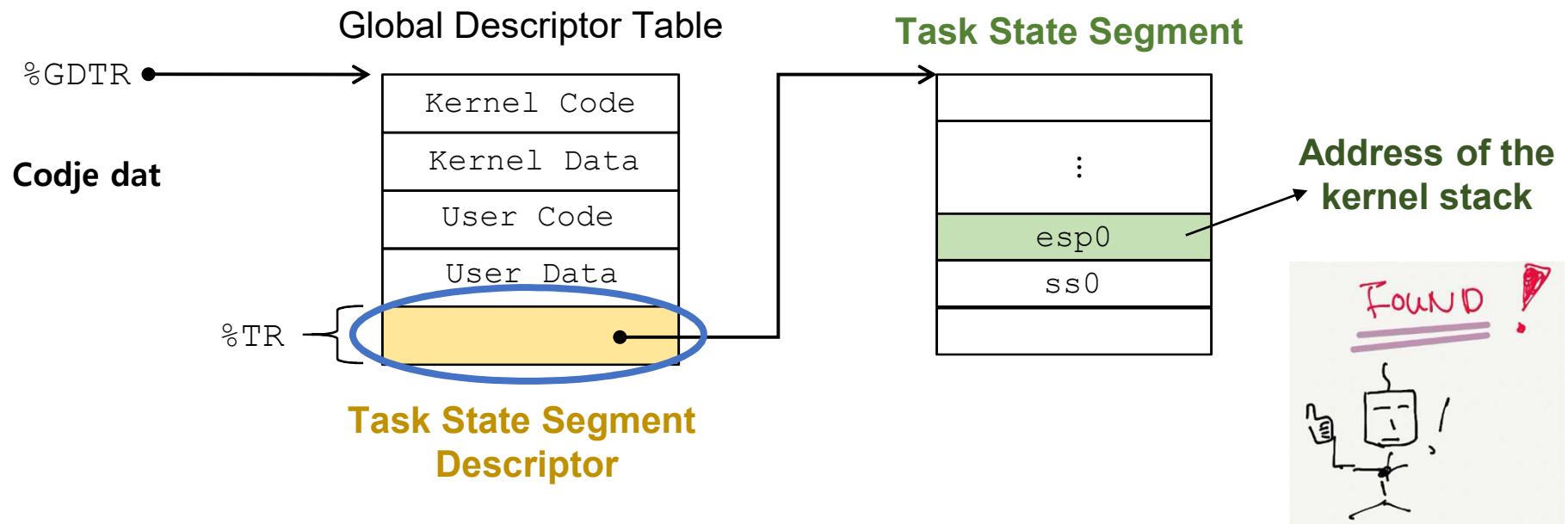
# Task State Segment (Cont.)

- Task State Segment has %esp & %ss register value of kernel mode process.
- esp0: Line 109
  - Address of the kernel stack
- ss0: Line 110
  - Segment Selector for stack.
  - Processor refers to it when push/pop instruction is executed.

```
106 // Task state segment format
107 struct taskstate {
108     uint link;           // Old ts selector
109     uint esp0;           // Stack pointers and segment selectors
110     ushort ss0;          // after an increase in privilege level
111     ushort padding1;
112
113     ...
145 };
```

# Finding the Task State Segment

- GDTR register & TS register.
  - %GDTR: the start address of global descriptor table.
  - %TR: the index of **task state segment descriptor** in the global descriptor table.
  - This descriptor points the **Task State Segment**.



# Setting the Task State Segment Descriptor

- switchuvm() switches the stack from kernel to user process p.
  - SEG\_TSS is 5.
  - It sets the entry at index 5 in GDT to Task state segment (Line 167 ~ 168).
  - It sets the value of TS register to 5 (Line 175).

```
157 void switchuvm(struct proc *p) {  
    ...  
166     pushcli();  
167     mycpu() ->gdt[SEG_TSS] = SEG16(STS_T32A, &mycpu() ->ts,  
168                               sizeof(mycpu() ->ts) - 1, 0);  
169     mycpu() ->gdt[SEG_TSS].s = 0;  
170     mycpu() ->ts.ss0 = SEG_KDATA << 3;  
171     mycpu() ->ts.esp0 = (uint)p->kstack + KSTACKSIZE;  
    ...  
175     ltr(SEG_TSS << 3);  
176     lcr3(V2P(p->pgdir)); // switch to process's address space  
177     popcli();  
178 }
```

# Setting the Task State Segment

- switchuvm() switches the stack from kernel to user.
  - It stores the value of current values of ss0 and esp0.
  - p->kstack stores the start address of kernel stack page.

```
157 void switchuvm(struct proc *p) {  
    ...  
166     pushcli();  
167     mycpu() ->gdt[SEG_TSS] = SEG16(STS_T32A, &mycpu() ->ts,  
168                               sizeof(mycpu() ->ts)-1, 0);  
169     mycpu() ->gdt[SEG_TSS].s = 0;  
170     mycpu() ->ts.ss0 = SEG_KDATA << 3;  
171     mycpu() ->ts.esp0 = (uint)p->kstack + KSTACKSIZE;  
    ...  
175     ltr(SEG_TSS << 3);  
176     lcr3(V2P(p->pgdir)); // switch to process's address space  
177     popcli();  
178 }
```

# Changing the stack

- CPL <= DPL:
  - The process can start this point only when the above condition is satisfied.
  - CPL == DPL:
    - The privilege level is not updated, stack change is not required.
  - CPL < DPL:
    - Save the current value of %ss & %esp, temporarily.
    - Fetch the new value to %ss & %esp from Task State Segment.
    - Push original value of %ss & %esp to new stack.

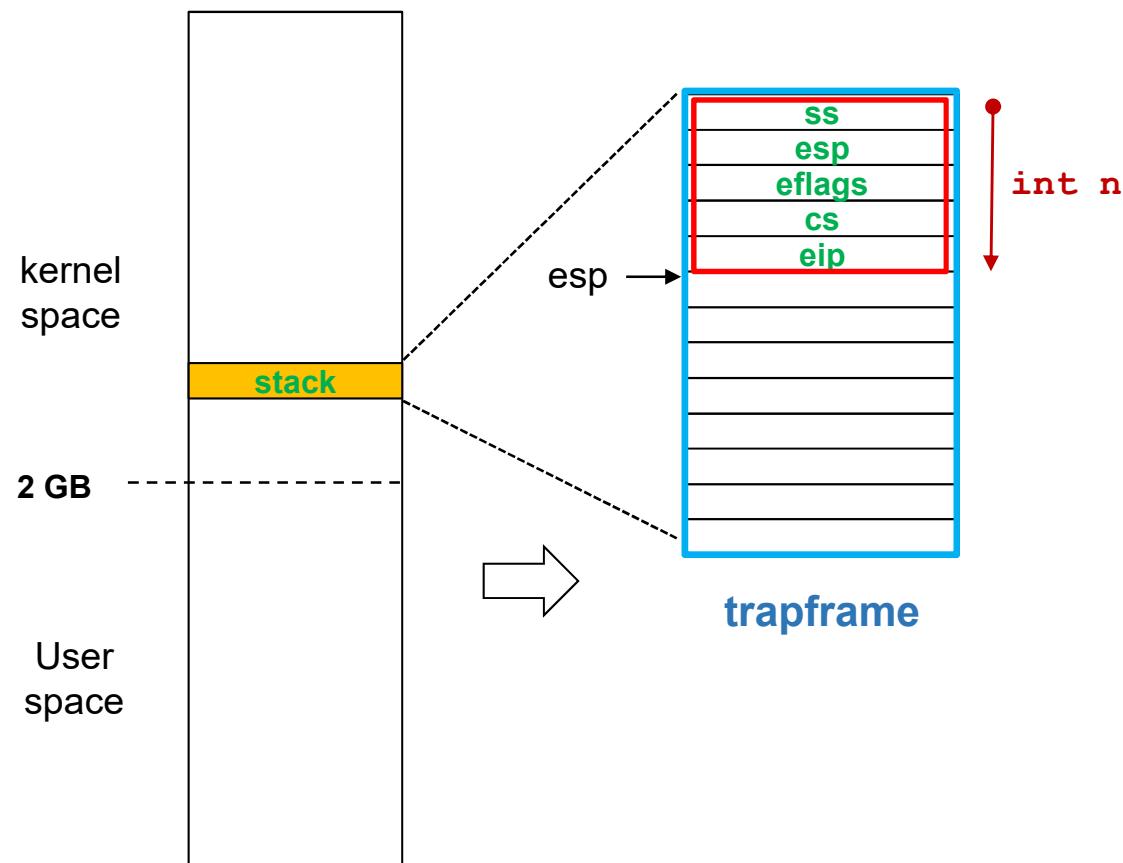
# int n

---

- ① Find the address of interrupt handler and check the privilege level.
- ② Change the stack from the user stack to the kernel stack.
- ③ Push the registers at the kernel stack.**
- ④ Load the address of interrupt handler to %eip register.

# Pushing some registers' value

- Push %eflags, %cs, and %eip to kernel stack.



## struct trapframe

- First, the values of five registers are pushed by hardware (%ss, %esp, %eflags, %cs, and %eip).
- The remaining values are pushed later...

```
150 struct trapframe {  
    ...  
171  
172     // below here defined by x86 hardware  
173     uint err;  
174     uint eip;  
175     ushort cs;  
176     ushort padding5;  
177     uint eflags;  
178  
179     // below here only when crossing rings,  
         such as from user to kernel  
180     uint esp;  
181     ushort ss;  
182     ushort padding6;  
183 };
```

Stack grows.

Push and popped by hardware

int n

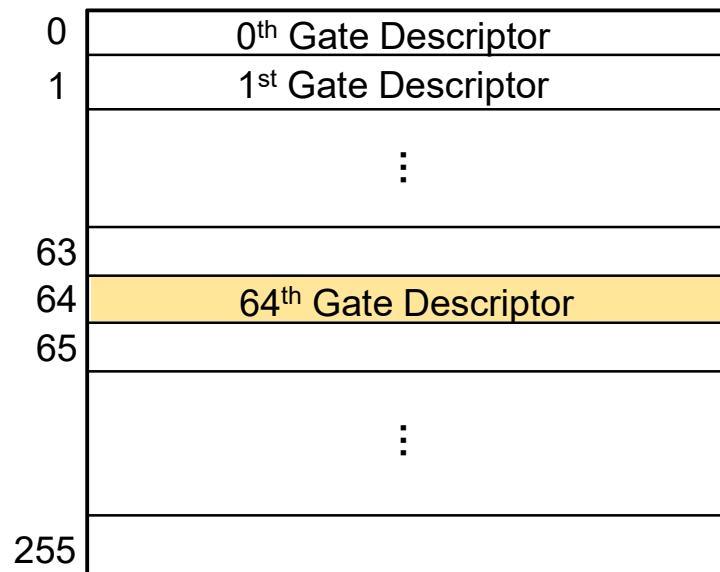
# int n

---

- ① Find the address of interrupt handler and check the privilege level.
- ② Change the stack from the user stack to the kernel stack.
- ③ Push some value of the registers.
- ④ Load the address of interrupt handler to %eip.

# Load the address of interrupt handler.

- Load %eip and %cs from the gate descriptor.
- 64<sup>th</sup> Gate Descriptor points to address of vector64.



System Call

%eip → vector64

%cs → SEG\_KCODE <<3

Interrupt Descriptor Table

```
24 SETGATE(idt[T_SYSCALL], 1, SEG_KCODE<<3, vectors[T_SYSCALL], DPL_USER);
```

at tvinit().

# int n

① Find the address of interrupt handler and check the privilege level.

- Fetch  $n^{\text{th}}$  gate descriptor from IDT.
- Examine if current privilege (in %cs) is lower than or equal to DPL of gate descriptor.

② Change the stack from the user stack to the kernel stack.

- If current privilege is lower than DPL,
  - load %ss & %esp from task state segment.
  - push original value of %ss and %esp.

③ Push some value of the registers.

- Push %eflags, %cs, and %eip to the stack.

④ Load the address of interrupt handler.

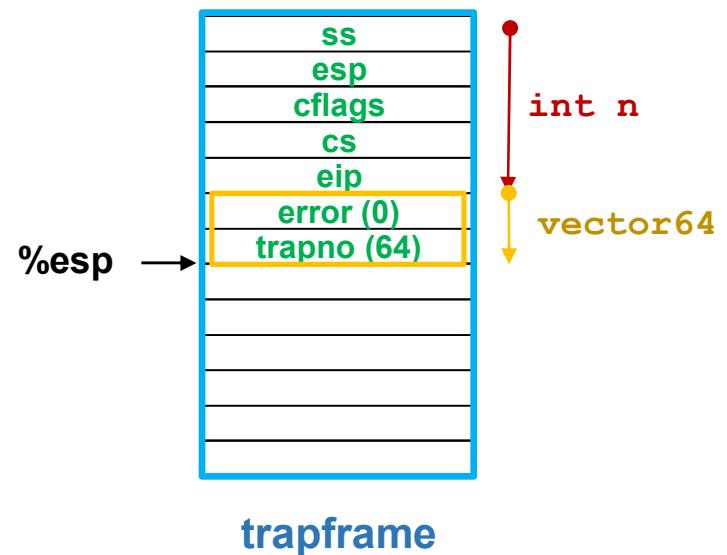
- %cs, %eip are loaded the  $n^{\text{th}}$  gate descriptor in IDT.
- Current privilege level is increased due to update of %cs.
- Begin to execute interrupt handler (vector64 for system call).

# Trap handler

- Pushes error number and trap number to the trap frame.
- Example: **vector 64 (system call)**
  - It pushes 0. This is value of error of the trapframe (Line 319).
  - It pushes 64. This is trap number of the trapframe (Line 320).
  - Then, call the function `alltraps`.

New value of %eip = address of **vector64**

```
# vectors.S
317 .globl vector64
318 vector64:
319     pushl $0
320     pushl $64
321     jmp alltraps
```



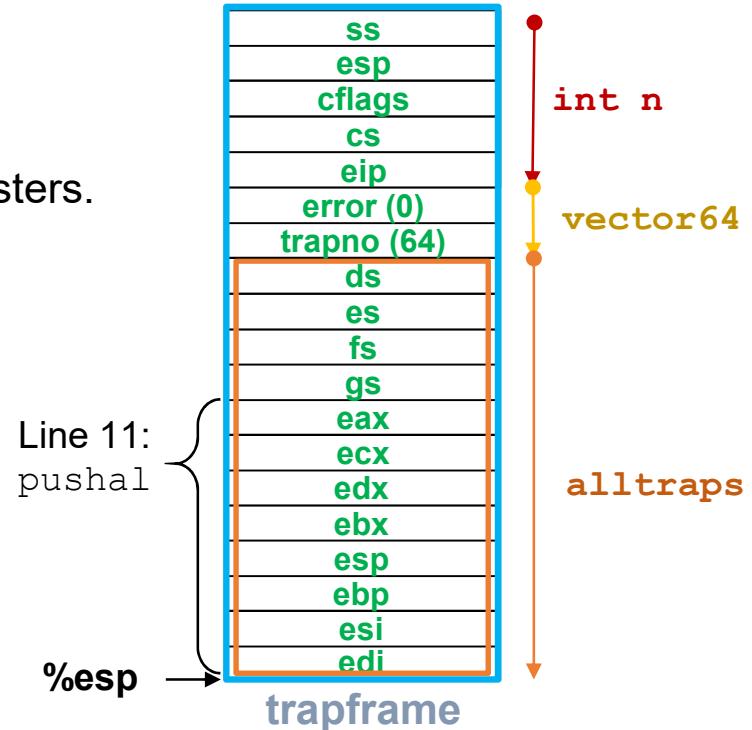
# Trap handler: alltraps

- Pushes the rest of the registers at the trap frame.
- Jump to trap handler.

```
4 .globl alltraps
5 alltraps:
6 # Build trap frame.
7 pushl %ds
8 pushl %es
9 pushl %fs
10 pushl %gs
11 pushal
12
13 # Set up data segments.
14 movw $(SEG_KDATA<<3), %ax
15 movw %ax, %ds
16 movw %ax, %es
17
18 # Call trap(tf), where tf=%esp
19 pushl %esp
20 call trap
21 addl $4, %esp
```

Push the value of registers.

call the function trap.



# Trap handler: alltraps (Cont.)

- Change the value of segment selector to access kernel's address space.

```
4 .globl alltraps
5 alltraps:
6 # Build trap frame.
7 pushl %ds
8 pushl %es
9 pushl %fs
10 pushl %gs
11 pushal
12
13 # Set up data segments.
14 movw $(SEG_KDATA<<3), %ax
15 movw %ax, %ds
16 movw %ax, %es
17
18 # Call trap(tf), where tf=%esp
19 pushl %esp
20 call trap
21 addl $4, %esp
```

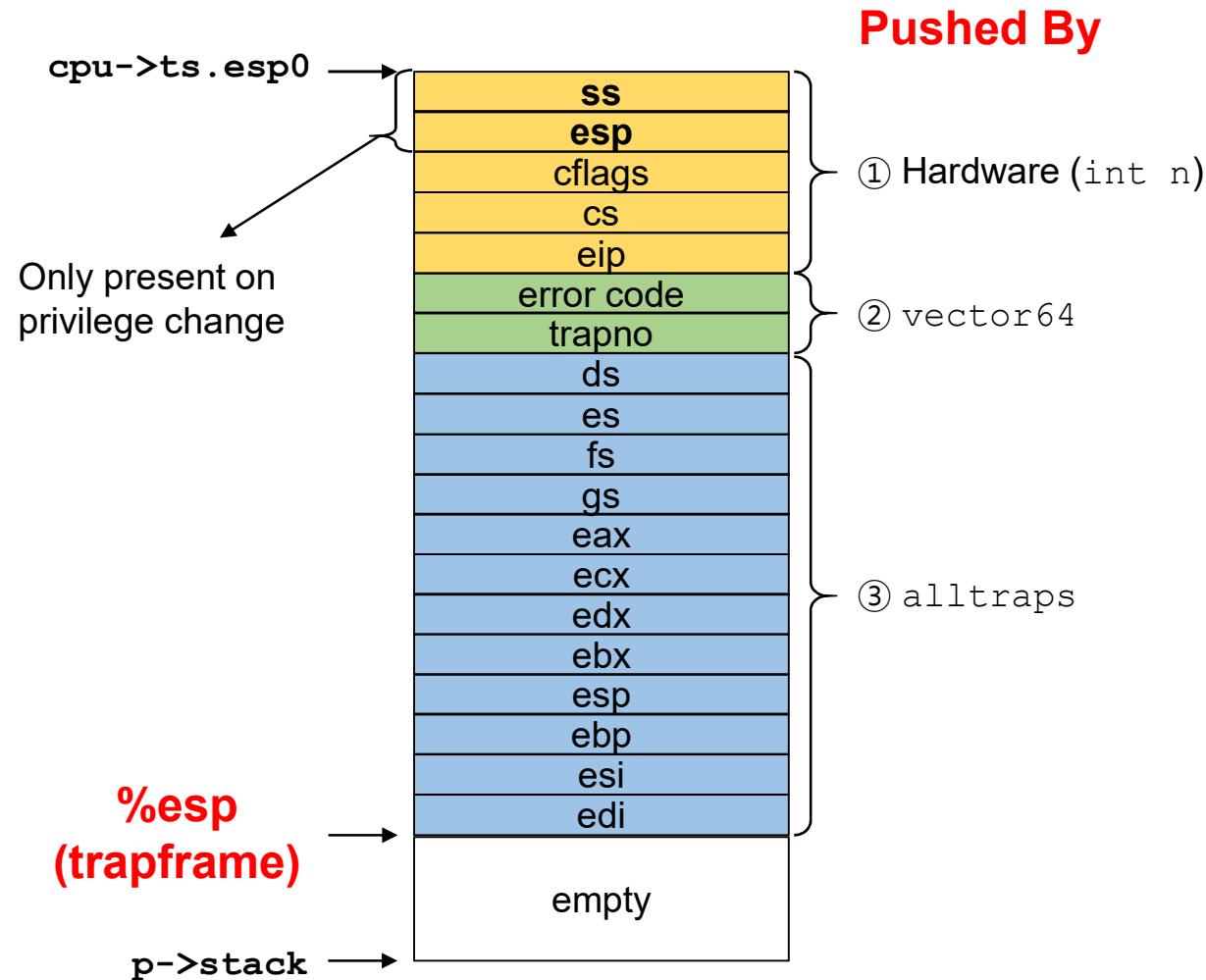
Change the value of segment selector

%ds & %es

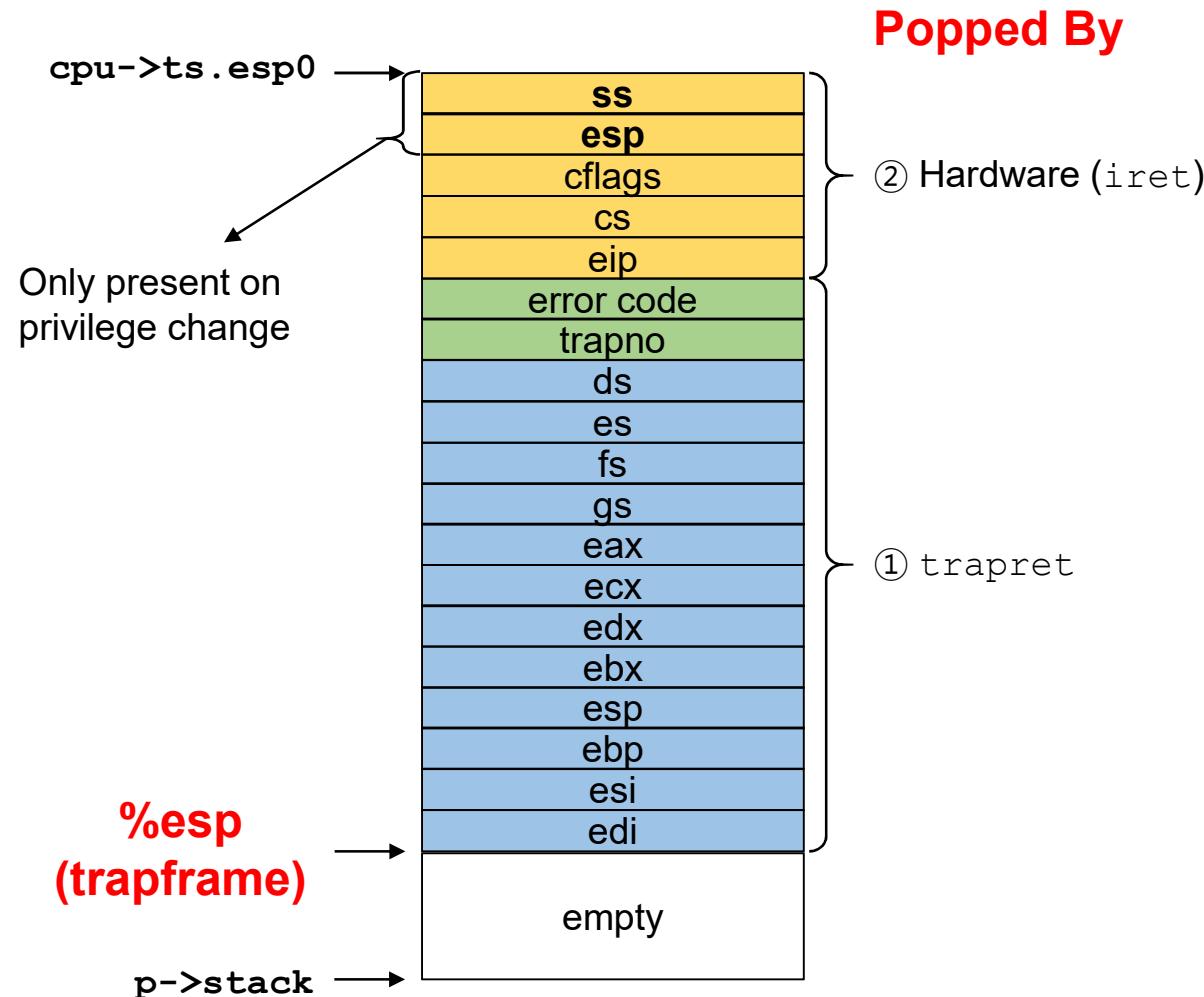
(SEG\_UDATA<<3 | DPL\_USER)

→ SEG\_KDATA<<3

# Kernel Stack Contents when entering the kernel



# Kernel Stack Contents when getting out of the kernel



## trap()

- `trameframe` is the parameter to `trap()` function.
- `%esp` points to the start (lowest) address of `trapframe`.
- pushes `%esp` at the stack before calling `trap()`.

```
18  # Call trap(tf), where tf=%esp
19  pushl %esp // Pointer to trapframe.
20  call trap // jump to 'real' handler.
21  addl $4, %esp
22
23  # Return falls through to trapret...
24 .globl trapret
25 trapret:
26  popal
27  popl %gs
28  popl %fs
29  popl %es
30  popl %ds
31  addl $0x8, %esp # trapno and errcode
32  iret
```

# trap ()

```
36 void
37 trap(struct trapframe *tf)
38 {
39     if(tf->trapno == T_SYSCALL) {
40         if(myproc()->killed)
41             exit();
42         myproc()->tf = tf;
43         syscall();
44         if(myproc()->killed)
45             exit();
46         return;
47     }
48
49     switch(tf->trapno) {
50     case T_IRQ0 + IRQ_TIMER:
51         ...
52     }
```

If trap number is T\_SYSCALL,  
trap() calls the syscall().

```
59     case T_IRQ0 + IRQ_IDE:
60         ...
61     case T_IRQ0 + IRQ_IDE+1:
62         ...
63     case T_IRQ0 + IRQ_KBD:
64         ...
65     case T_IRQ0 + IRQ_COM1:
66         ...
67     case T_IRQ0 + 7:
68     case T_IRQ0 + IRQ_SPURIOUS:
69         ...
70     default:
71         ...
72     }
```

# Code: The first system call exec()

- initcode.S
  - Process pushes the arguments for an exec call on the process's stack.
  - Put the system call number in %eax
  - The system call numbers match the entries in the syscalls array, a table of function pointers.

```
9 .globl start
10 start:
11     pushl $argv          Fake return address
12     pushl $init
13     pushl $0    // where caller pc would be
14     movl $SYS_exec, %eax
15     int $T_SYSCALL
16
```

```
107 static int (*syscalls[])(void) = {
108 [SYS_fork]      sys_fork,
109 [SYS_exit]      sys_exit,
110 [SYS_wait]      sys_wait,
111 [SYS_pipe]      sys_pipe,
112 [SYS_read]      sys_read,
113 [SYS_kill]      sys_kill,
114 [SYS_exec]      sys_exec, // Line 114
115 [SYS_fstat]     sys_fstat,
116 [SYS_chdir]     sys_chdir,
117 [SYS_dup]       sys_dup,
118 [SYS_getpid]   sys_getpid,
119 [SYS_sbrk]      sys_sbrk,
120 [SYS_sleep]    sys_sleep,
121 [SYS_uptime]   sys_uptime,
122 [SYS_open]      sys_open,
123 [SYS_write]    sys_write,
124 [SYS_mknod]    sys_mknod,
125 [SYS_unlink]   sys_unlink,
126 [SYS_link]     sys_link,
127 [SYS_mkdir]    sys_mkdir,
128 [SYS_close]    sys_close,
129};
```

syscall.c

# syscall()

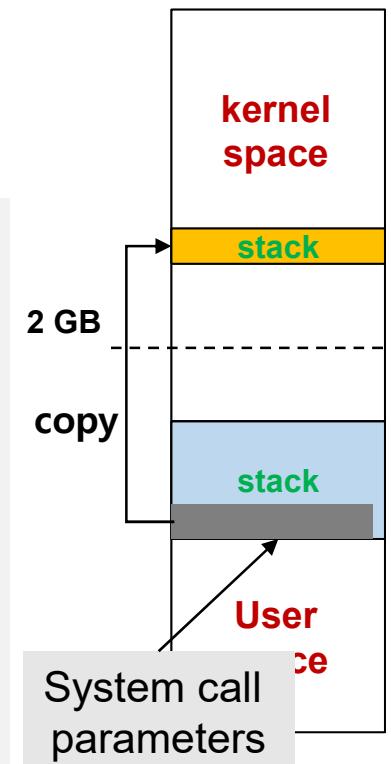
- checks the system call number at %eax.
- Then executes the proper function.

```
131 void
132 syscall(void)
133 {
134     int num;
135     struct proc *curproc = myproc();
136
137     num = curproc->tf->eax;
138     if(num > 0 && num < NELEM(syscalls) && syscalls[num]) {
139         curproc->tf->eax = syscalls[num]();
140     } else {
141         cprintf("%d %s: unknown sys call %d\n",
142                 curproc->pid, curproc->name, num);
143         curproc->tf->eax = -1;
144     }
145 }
```

# The parameter to system call

- Parameters to the system call are stored in the caller's stack.
- If the caller is user, the parameters are initially in the user stack.
- It copies the parameter to system call to its local variable.
  - argint(): copy the integer value to local variable.
  - argptr(): copy the pointer value to local variable.
  - argstr(): copy the string to local variable.

```
396 int
397 sys_exec(void)
398 {
399     char *path, *argv[MAXARG];
400     int i;
401     uint uargv, uarg;
402
403     if(argstr(0, &path) < 0 || argint(1, (int*)&uargv) < 0) {
404         return -1;
405     }
...
}
```



# Accessing caller's stack

- System call accesses the caller's stack by using trapframe.
- The trapframe has caller's stack address at the `esp` variable.

```
49 int
50 argint(int n, int *ip)
51 {
52     return fetchint((myproc() ->tf ->esp) + 4 + 4*n, ip);
53 }

...
56 int
57 argstr(int n, char **pp)
58 {
59     int addr;
60     if(argint(n, &addr) < 0)
61         return -1;
62     return fetchstr(addr, pp);
63 }
```

# Summary

---

- `int instruction`

The current privilege level should be lower than or equal to DPL of gate descriptor of IDT.

The `%esp` & `%ss` are loaded from task state segment pointed by segment descriptor in GDT.

The `%eip` & `%cs` are loaded from gate descriptor of IDT.

- `alltraps` & `trap()`

`alltraps` pushes the value of registers and calls `trap()`.

`trap()` calls the function with respect to trap number.

If trap number is 64 (system call), it calls `syscall()`.

- `syscall()`

`syscall()` calls the proper function with respect to the value of `%eax`.