# Chapter 5 Kernel Code (kernel)

## 5.1 Introduction

Directory linux/kernel/ contains 10 C language file and 2 assembly language file and a Makefile for manage compiling as listed in List 5-1. The files in three subdirectories will be studied in the following chapters.

This chapter mainly comments about these 13 files. Firstly, we summaris the basic functions for each source code file to gain a general understanding of them, then we study each file one by one.

List 5-1 files in linux/kernel/ directory

| | File name | Size | Last modified time(GMT) | Description |
|---|---|---|---|---|
| | blk_drv/ | | 1991-12-08 14:09:29 | |
| | chr_drv/ | | 1991-12-08 18:36:09 | |
| | math/ | | 1991-12-08 14:09:58 | |
| | Makefile | 3309 bytes | 1991-12-02 03:21:37 | m |
| | asm.s | 2335 bytes | 1991-11-18 00:30:28 | m |
| | exit.c | 4175 bytes | 1991-12-07 15:47:55 | m |
| | fork.c | 3693 bytes | 1991-11-25 15:11:09 | m |
| | mktime.c | 1461 bytes | 1991-10-02 14:16:29 | m |
| | panic.c | 448 bytes | 1991-10-17 14:22:02 | m |
| | printk.c | 734 bytes | 1991-10-02 14:16:29 | m |
| | sched.c | 8242 bytes | 1991-12-04 19:55:28 | m |
| | signal.c | 2651 bytes | 1991-12-07 15:47:55 | m |
| | sys.c | 3706 bytes | 1991-11-25 19:31:13 | m |
| | system_call.s | 5265 bytes | 1991-12-04 13:56:34 | m |
| | traps.c | 4951 bytes | 1991-10-30 20:20:40 | m |
| | vsprintf.c | 4800 bytes | 1991-10-02 14:16:29 | m |

## 5.2 Main Functions

The source code in this directory can be seperated into three classes. The first one contains files used to deal with interrupts generated by CPU execption, hardware request and software instructions. The second class is the files for system call services. The last class contains process generating, scheduling, exiting and some other assistant code files. The relationship of these files is demonstrated in Figure 5-1.
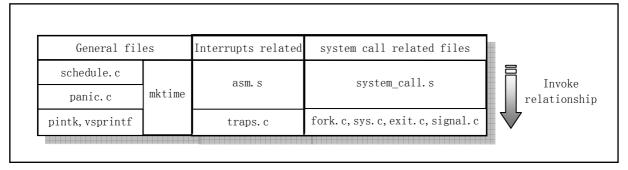
| General files | | Interrupts related | system call related files |
|---|---|---|---|
| schedule.c | mktime | asm.s | system_call.s |
| panic.c | | | |
| pintk, vsprintf | | traps.c | fork.c, sys.c, exit.c, signal.c |

Invoke relationship

Figure 5-1 The relationship of files in kernel/ directory

## 5.2.1 Interrupt processing

There are two files related to interrupts handling: asm.s and traps.c. asm.s includes most part of assembly codes for hardware execption handling, while traps.c includes most C functions that called from the interrupt assembly codes. This C functions sometimes called Bottom-halfs of the interrupt handling part. In addition, several other interrupt handling functions are implemented in system_call.s and mm/page.s file.

An interrupt gate or trap gate points indirectly to a procedure which will execute in the context of the currently executing task. The selector of the gate points to an executable-segment descriptor in either the GDT or the current LDT. The offset field of the gate points to the beginning of the interrupt or exception handling procedure. Just as with a control transfer due to a CALL instruction, a control transfer to an interrupt or exception handling procedure uses the stack to store the information needed for returning to the original procedure.

Before the CPU transfer the control to the interrupt handling procedure when a interrupt occurred, it will first push at least 12 bytes onto the stack. The other similar way as a long CALL between segments is that the CPU pushes the information onto stack of target code instead of the stack of being interrupted code. Thus, when a interrupt occurred, Linux uses kernel mode stack. In addition, an interrupt always pushes the EFLAGS register onto the stack. If there is privilege transition, for example from user to kernel privilege, CPU will also push the orignal stack pointer used by the currently executing task onto the target stack as shown in Figure 5-2.

| Old SS |
|---|
| Old ESP |
| Old EFLAGS |
| Old CS |
| Old EIP |

| Old SS |
|---|
| Old ESP |
| Old EFLAGS |
| Old CS |
| Old EIP |
| Error Code |

push direction

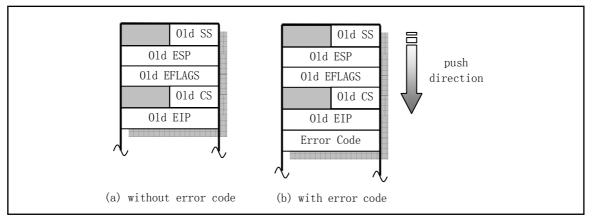(a) without error code    (b) with error code

Figure 5-2 Stack layout after exception of interrupt with privilege transition

asm.s program is used to handle Intel reserved interrupts 0x00-0x1f. For Linux 0.11, the programmerable interrupt controllers (PIC) were initialized to use interrupt vectors with in range of 0x20-0x2f and the corresponding interrupt handling procedures are scattered in each initialization routines of hardware drivers. The handling of Linux system call interrupt int 0x80 is in file kernel/system_call.s. For the predetermined identifiers, see the table after the program listing.

As some software exceptions that relate to a specific segment will cause CPU to generate a error code internally, program asm.s clarify the processing of all interrupts into two classes based on whether a interrupt contains a error code. But the handling procedure of them are similar as depictured in Figure 5-3.
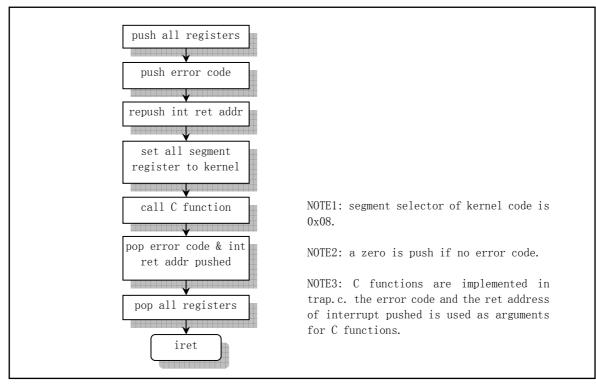


Figure 5-3 The flow chart of hardware exception (fault, traps) handling

## 5.2.2  System calls handling

In Linux system, application program that use the function of kernel is through the system call, that is the interrupt 0x80. The register eax contains the call function number and ebx, ecx and edx were used to store arguments if there is any. In Linux 0.11 kernel, there are 72 system calls implemented in programs system_call.s, fork.c, signal.c, sys.c and exit.c.

The way system_call.s program handling system call interrupt is similar as in asm.s program. In addition, it also contains the handling procedure of timer, hard disk and floppy drives. fork.c and signal.c are used to privide C functions for corresponding interrupt handling just like the operation of traps.c. The primary functions in fork.c is find_empty_process() and copy_process(). For the signal.c, see the discussing at the beginning of its program listing.

sys.c and exit.c program implements numerous other system calls C functions (sys_xxx()). We can properly know the using of a particular functions by its naming method. For example,

a C function name begins with 'do_' is usually a general function invoked in the process of a system call; The C function begins with 'sys_' is usually only called in a particular system call interrupt. For example, do_signal() is invoked by almost all system calls in system_call.s, while sys_execve() is only invoked in a particular system call.

### 5.2.3  Other programs

These programs include schedule.c, mktime.c, panic.c printk.c and vsprintf.c.

schedule.c contains the most important routines: schedule(), sleep_on() and wakeup() used to switch and change the status of tasks. In addition, there are several routines related to timer handling for timer interrupt and floppy drives. mktime.c program is used to calculate the start time of kernel and only invoked in init/main.c. panic.c contains a routine panic() and used to display vital error occurred in kernel code and stop the kernel. printk.c and vsprintf.c are support programs used to display messages of the kernel.

## 5.3  Makefile file

### 5.3.1  Functions

kernel/Makefile is a configuration file used by make tools for compiling the programs in kernel/ directory, not including the programs in three subdirectories. This file is almost identical with the file listed in program 2-1 in chapter 2.

### 5.3.2  Comments

Program 5-1 linux/kernel/Makefile

```
 1 #
 2 # Makefile for the FREAX-kernel.
 3 #
 4 # Note! Dependencies are done automagically by 'make dep', which also
 5 # removes any old dependencies. DON'T put your own dependencies here
 6 # unless it's something special (ie not a .c file).
 7 #
   # (The Initial name of Linux is called FREAX an was changed to 'Linux' by the administrator
   # of ftp.funet.fi)
 8
 9 AR      =gar    # GNU ar program used to create, modify, and extract files from archives.
10 AS      =gas    # GNU assembler.
11 LD      =gld    # GNU linker used to combine a number of object and archive files.
12 LDFLAGS =-s -x  # Options for gld. -s strip all symbols; -x delete all local symbols.
13 CC      =gcc    # GNU C language compiler.
14 CFLAGS  =-Wall -O -fstrength-reduce -fomit-frame-pointer -fcombine-regs \
15          -finline-functions -mstring-insns -nostdinc -I../include
   # Options for gcc.
   # -Wall show all warning messages;
   # -O optimze code length and executing time;
   # -fstrength-reduce optimzes codes in loop, eliminates iteration variables;
   # -fomit-frame-pointer Don't keep frame pointer in a register for functions that don't need;
   # -fcombine-regs Combine registers, minimum the using of register classes;
```

```
   # -mstring-insns This option is added by Linus to gcc 1.40 used to adopt string instructions
   # in i386.
   # -nostdinc Don't use header files in default directory, e.g. /usr/include/
   # -I../include Indicates the directory of header files used (../include).
16 CPP    =gcc -E -nostdinc -I../include
   # Options for C preprocessor.
   # -E only runs the C preprocessing, The output is in the form of preprocessed source code, which
   # is sent to the standard output;
17
   # The following role (L18) is used to indicate make program to use the command (L19) below
   # to compile all '.c' files into '.s' assembly files.
   # The command (L19) indicates gcc to compile C code by using the options provided by 'CFLAGS',
   # but don't run linking process (-S). Thus producing assembler files coresponding each '.c'
   # files with suffix of '.s'.
   # -o indicats the following is name of output; '$*.s' (or '$@') is automatic variable.
   # '$<' represents the first precondition, e.g. the '*.c' files here.
18 .c.s:
19         $(CC) $(CFLAGS) \
20         -S -o $*.s $<
   # The following role will compile all '.s' assembler file into '.o' object files using gas.
21 .s.o:
22         $(AS) -c -o $*.o $<
23 .c.o:                        # '*.c' files to '*.o' object files, don't do linking.
24         $(CC) $(CFLAGS) \
25         -c -o $*.o $<
26
27 OBJS  = sched.o system_call.o traps.o asm.o fork.o \      # Define Object variables OBJS.
28        panic.o printk.o vsprintf.o sys.o exit.o \
29        signal.o mktime.o
30
   # Linking to produce target object 'kernel.o' after precondition 'OBJS' is satisfied.
31 kernel.o: $(OBJS)
32         $(LD) -r -o kernel.o $(OBJS)
33         sync
34
   # The following role is used to do cleaning after once compile. When we execute 'make clean'
   # at shell's command prompt, the 'make' will run the commands below (L36-40), erasing all
   # files generated by compiler.
   # 'rm' is file deletion command. Option '-f' neglecting non-exist files and errors.
35 clean:
36         rm -f core *.o *.a tmp_make keyboard.s
37         for i in *.c;do rm -f `basename $$i .c`.s;done
38         (cd chr_drv; make clean)
39         (cd blk_drv; make clean)
40         (cd math; make clean)
41
   # The following role is used to check the dependence between each source files:
   # Use stream editor 'sed' to deal with 'Makefile' (this file), the output is a temporary
   # file 'tmp_make', in which, all lines after '### Dependencies' (begins from L51) were
   # deleted. Then, for all '.c' files in the directory, output a file name, which suffix is
   # converted to '.s' and followed by a white space character, and, afterwards, executing
   # preprocess programm (CPP) to produce a dependent line for the corresponding '.c' file as
   # illustrated in line L51-55 for exit.c file. Finally, copy the temporary file to 'Makefile'
```

```
     # to renew it.
     # Option '-M' tells the CPP to output the roles for each relating object file in make syntax.
42 dep:
43          sed '/\#\#\# Dependencies/q' < Makefile > tmp_make
44          (for i in *.c;do echo -n `echo $$i | sed 's,\.c,\.s,'`" "; \
45                $(CPP) -M $$i;done) >> tmp_make
46          cp tmp_make Makefile
47          (cd chr_drv; make dep)
48          (cd blk_drv; make dep)
49
50 ### Dependencies:
51 exit.s exit.o : exit.c ../include/errno.h ../include/signal.h \
52   ../include/sys/types.h ../include/sys/wait.h ../include/linux/sched.h \
53   ../include/linux/head.h ../include/linux/fs.h ../include/linux/mm.h \
54   ../include/linux/kernel.h ../include/linux/tty.h ../include/termios.h \
55   ../include/asm/segment.h
56 fork.s fork.o : fork.c ../include/errno.h ../include/linux/sched.h \
57   ../include/linux/head.h ../include/linux/fs.h ../include/sys/types.h \
58   ../include/linux/mm.h ../include/signal.h ../include/linux/kernel.h \
59   ../include/asm/segment.h ../include/asm/system.h
60 mktime.s mktime.o : mktime.c ../include/time.h
61 panic.s panic.o : panic.c ../include/linux/kernel.h ../include/linux/sched.h \
62   ../include/linux/head.h ../include/linux/fs.h ../include/sys/types.h \
63   ../include/linux/mm.h ../include/signal.h
64 printk.s printk.o : printk.c ../include/stdarg.h ../include/stddef.h \
65   ../include/linux/kernel.h
66 sched.s sched.o : sched.c ../include/linux/sched.h ../include/linux/head.h \
67   ../include/linux/fs.h ../include/sys/types.h ../include/linux/mm.h \
68   ../include/signal.h ../include/linux/kernel.h ../include/linux/sys.h \
69   ../include/linux/fdreg.h ../include/asm/system.h ../include/asm/io.h \
70   ../include/asm/segment.h
71 signal.s signal.o : signal.c ../include/linux/sched.h ../include/linux/head.h \
72   ../include/linux/fs.h ../include/sys/types.h ../include/linux/mm.h \
73   ../include/signal.h ../include/linux/kernel.h ../include/asm/segment.h
74 sys.s sys.o : sys.c ../include/errno.h ../include/linux/sched.h \
75   ../include/linux/head.h ../include/linux/fs.h ../include/sys/types.h \
76   ../include/linux/mm.h ../include/signal.h ../include/linux/tty.h \
77   ../include/termios.h ../include/linux/kernel.h ../include/asm/segment.h \
78   ../include/sys/times.h ../include/sys/utsname.h
79 traps.s traps.o : traps.c ../include/string.h ../include/linux/head.h \
80   ../include/linux/sched.h ../include/linux/fs.h ../include/sys/types.h \
81   ../include/linux/mm.h ../include/signal.h ../include/linux/kernel.h \
82   ../include/asm/system.h ../include/asm/segment.h ../include/asm/io.h
83 vsprintf.s vsprintf.o : vsprintf.c ../include/stdarg.h ../include/string.h
```

# 5.4  asm.s

## 5.4.1  Functions

asm.s program involves the low level codes for processing exceptions detected by CPU, includes the handling of FPU exceptions. This program is intimately related to kernel/traps.c

program. Its primary handling method is to invoke the corresponding C functions to do the detail processing, display the error location and error code, and finally exit.

When start reading this piece of code, it is helpful to refer to the changes in stack as illustrated in Figure 5-4. In this figure, each line of frames represents 4 bytes. Before executes the code, the stack pointer esp points to the interrupt return address (esp0 in figure). After it pushes the address of C function do_divide_error() or other related C funtions into the stack, stack pointer moves to the location esp1 in figure. At this moment, the function's address is swapped into register eax while the value of eax is swapped into stack. After several other registers pushed into stack, the location of stack pointer is at esp2. Just before invoking the do_divide_error(), it pushes the esp0 into stack as a argument for C function, that is at the location of esp3 in figure. When returns from the C function, it adds value 8 to the current stack pointer to back to the location of esp2.
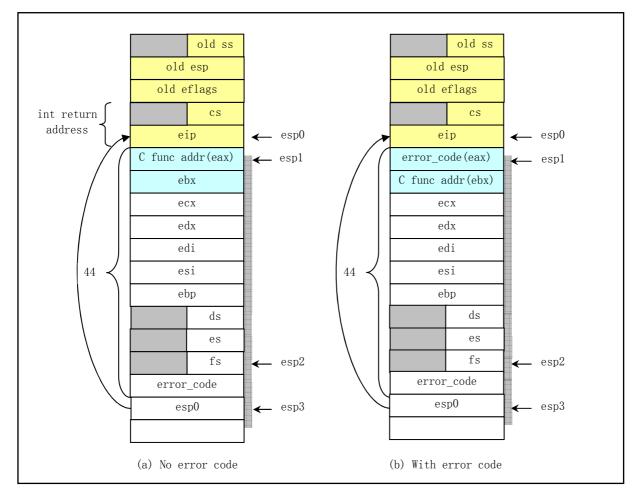


Figure 5-4 Stack map of error handling for exceptions

As stated above, the error code and the value of stack pointer when just enter the program (esp0) are used as arguments for the C function (e.g. do_divide_error()). In the program traps.c, the prototype of do_divide_error() is:

void do_divide_error(long esp, long error_code)

Therefore, this function can print out the error location and error code of the program causing exception.

## 5.4.2 Comments

Program 5-2 linux/kernel/asm.s

```
1  /*
2   *  linux/kernel/asm.s
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  /*
8   * asm.s contains the low-level code for most hardware faults.
9   * page_exception is handled by the mm, so that isn't here. This
10  * file also handles (hopefully) fpu-exceptions due to TS-bit, as
11  * the fpu must be properly saved/resored. This hasn't been tested.
12  */
13
   # This program is used to handling the Intel reserved interrupts 0x00-0x16.
   # The following lines are declarations of global functions, implemented in traps.c.
14 .globl _divide_error,_debug,_nmi,_int3,_overflow,_bounds,_invalid_op
15 .globl _double_fault,_coprocessor_segment_overrun
16 .globl _invalid_TSS,_segment_not_present,_stack_segment
17 .globl _general_protection,_coprocessor_error,_irq13,_reserved
18
   # The following codes were used to handle exceptions with no error code. Refers to figure 5-4(a).
   # The label '_do_divide_error' is actually a internal name representing the coressponding C
   # function do_divide_error().
   #
   # int 0 - Procedure for handling devide by zero error. (Fault)
19 _divide_error:
20        pushl $_do_divide_error    # Push the addr of C function onto stack first.
                                     # AT&T immediate operands are preceded by '$'.
21 no_error_code:                    # Entry point for interrupts w/o error code, Refer to L55.
22        xchgl %eax,(%esp)          # Exchange the contents of stack pointed by esp and eax.
23        pushl %ebx                 # That is: addr of do_divide_error() ->eax, eax ->(esp)
24        pushl %ecx
25        pushl %edx
26        pushl %edi
27        pushl %esi
28        pushl %ebp
29        push %ds                   # 16-bit segment register will occupy 4 bytes on stack.
30        push %es
31        push %fs
   # Pushes error code and the esp of interrupt return addr as arguments for C function at L39.
32        pushl $0                   # "error code". push 0 to resemble a error code.
33        lea 44(%esp),%edx          # Push the stack pointer of interrupt return address.
34        pushl %edx
35        movl $0x10,%edx            # 0x10 is the value of kernel data segment selector.
36        mov %dx,%ds                # Prepare segment register to point to kernel
37        mov %dx,%es
38        mov %dx,%fs
```

```
39        call *%eax                  # Invokes the C function (e.g. do_divide_error()).
                                      # Note that AT&T absolute (as opposed to PC relative)
                                      # jump/call operands are prefixed by '*'
40        addl $8,%esp                # Get rid of two values on stack.
41        pop %fs
42        pop %es
43        pop %ds
44        popl %ebp
45        popl %esi
46        popl %edi
47        popl %edx
48        popl %ecx
49        popl %ebx
50        popl %eax                   # Pop the original eax.
51        iret
52
   # int1 -- Entry point for debug interrupt. (Fault/Trap)
   # Raised when the T flag of eflags is set or when the address of an instruction or operand
   # falls within the range of an active debug register
53 _debug:
54        pushl $_do_int3        # Push the address of do_debug().
55        jmp no_error_code
56
   # int2 -- Entry point for nonmaskable interrupts (NMI). (Trap)
57 _nmi:
58        pushl $_do_nmi
59        jmp no_error_code
60
   # int3 -- Entry point for breakpoint interrupt. (Trap)
   # Caused by an int3 instruction.Usually inserted by a debugger.
61 _int3:
62        pushl $_do_int3
63        jmp no_error_code
64
   # int4 -- Entry point for overflow issued interrupt. (Trap)
   # Raised when the OF (overflow) flag of eflags is set.
65 _overflow:
66        pushl $_do_overflow
67        jmp no_error_code
68
   # int5 -- Entry point for bounds check interrupt. (Fault)
   # A bound instruction is executed with the operand outside of the valid address bounds.
69 _bounds:
70        pushl $_do_bounds
71        jmp no_error_code
72
   # int6 -- Entry point for invalid opcode interrupt. (Fault)
   # The CPU execution unit has detected an invalid opcode.
73 _invalid_op:
74        pushl $_do_invalid_op
75        jmp no_error_code
76
   # int9 -- Entry point for coprocessor error. (Fault)
```

```
      # Problems with the external mathematical coprocessor.
77 _coprocessor_segment_overrun:
78         pushl $_do_coprocessor_segment_overrun
79         jmp no_error_code
80
      # int15 - Entry point for interrupt reserved by Intel.
81 _reserved:
82         pushl $_do_reserved
83         jmp no_error_code
84
      # int45 -- (= 0x20 + 13) Coprocessor generated interrupt, configured by Linux in PIC.
      # When copprocessor finishes an operation, it will issue IRQ 13 to tell CPU the state.
      # During the operation of 80387, the CPU will wait for its finish.
85 _irq13:
86         pushl %eax
87         xorb %al,%al
      # By write to I/O port 0xf0, this interrupt handler will remove the CPU's extension of BUSY
      # signal and reactivate the 80387's pin of coprocessor extended request (PEREQ).
      # The use of an interrupt for a copprocessor is to guarantee that an exception from a
      # coprocessor instruction will be detected and INTR 13 will be serviced prior to the
      # execution of any further 80387 instructions.
88         outb %al,$0xF0
89         movb $0x20,%al
90         outb %al,$0x20          # Send End of Interrupt (EOI) signal to 8259A master chip.
91         jmp 1f                  # Delay a while.
92 1:      jmp 1f
93 1:      outb %al,$0xA0          # Send EOI signal to 8259A slave chip.
94         popl %eax
95         jmp _coprocessor_error  # _coprocessor_error is at L131 in kernel/system_call.s.
96
      # When processing the following interrupts, CPU will push the error code onto stack after the
      # the interrupt return address. Therefore, we should pop up the error code when return.
      # Refers to figure 5-4(b).
      # int8 -- Entry point for double fault interrupt. (Abort)
      # Normally, when the CPU detects an exception while trying to call the handler for a prior
      # exception, the two exceptions can be handled serially. In a few cases, however, the processor
      # cannot handle them serially, so it raises this exception.
97 _double_fault:
98         pushl $_do_double_fault     # Push the pointor of C function onto stack.
99 error_code:
100        xchgl %eax,4(%esp)          # error code <-> %eax, old eax is saved onto stack.
101        xchgl %ebx,(%esp)           # &function <-> %ebx, old ebx is saved onto stack.
102        pushl %ecx
103        pushl %edx
104        pushl %edi
105        pushl %esi
106        pushl %ebp
107        push %ds
108        push %es
109        push %fs
110        pushl %eax                  # error code
111        lea 44(%esp),%eax           # offset
112        pushl %eax
```

```
113          movl $0x10,%eax                  # Setup kernel data segment register.
114          mov %ax,%ds
115          mov %ax,%es
116          mov %ax,%fs
117          call *%ebx                       # Invoke the C function. Arguments are on the stack.
118          addl $8,%esp                     # Get gid of arguments on the stack.
119          pop %fs
120          pop %es
121          pop %ds
122          popl %ebp
123          popl %esi
124          popl %edi
125          popl %edx
126          popl %ecx
127          popl %ebx
128          popl %eax
129          iret
130
     # int10 -- Entry point for invlid task status segment (TSS). (Fault)
     # The CPU has attempted a context switch to a process having an invalid Task State Segment.
131 _invalid_TSS:
132          pushl $_do_invalid_TSS
133          jmp error_code
134
     # int11 -- Entry point for segment not present interrupt. (Fault)
     # A reference was made to a segment not present in memory.
135 _segment_not_present:
136          pushl $_do_segment_not_present
137          jmp error_code
138
     # int12 -- Entry point for stack segment error interrupt. (Fault)
     # The instruction attempted to exceed the stack segment limit, or the segment identified by
     # ss is not present in memory.
139 _stack_segment:
140          pushl $_do_stack_segment
141          jmp error_code
142
     # int13 -- Entry point for general protection interrupt. (Fault).
     # One of the protection rules in the protected mode of the 80386 has been violated.
143 _general_protection:
144          pushl $_do_general_protection
145          jmp error_code
146
     # int7 -- _device_not_available is at L148 in kernel/system_call.s file.
     # int14 -- _page_fault) is at L14 in mm/page.s file.
     # int16 -- _coprocessor_error is at L131 in kernel/system_call.s file.
     # int 0x20 -- _timer_interrupt is at L176 in kernel/system_call.s file.
     # int 0x80 -- _system_call is at L80 in kernel/system_call.s file.
```

## 5.4.3 Information

### 5.4.3.1 Intel reserved interrupts

Table 5-1 presents the difinition of interrupt vectors reserved by Intel.

Table 5-1 Interrupt vectors reserved by Intel

| Int No. | Name | Type | Signal | Description |
|---|---|---|---|---|
| 0 | Devide error | Fault | SIGFPE | Occurred during a divide instruction when the divisor is zero. |
| 1 | Debug | Fault/ Trap | SIGTRAP | Raised when the T flag of eflags is set or when the address of an instruction or operand falls within the range of an active debug register. |
| 2 | nmi | Hardware | | Used for nonmaskable interrupts (NMI) |
| 3 | Breakpoint | Trap | SIGTRAP | The int 3 instruction causes this trap. |
| 4 | Overflow | Trap | SIGSEGV | This trap occurs when the processor encounters an INTO instruction and the OF (overflow) flag is set. |
| 5 | Bounds check | Fault | SIGSEGV | This fault occurs when the processor finds that the operand exceeds the specified limits while executing a bounds instruction. |
| 6 | Invalid Opcode | Fault | SIGILL | This fault occurs when an invalid opcode is detected by the CPU execution unit. |
| 7 | Device not available | Fault | SIGSEGV | This exception occurs in either of two conditions: (a) The CPU encounters an ESC instruction, and the EM (emulate) bit of CR0 is set; (b) The CPU encounters either the WAIT instruction or an ESC instruction, and both the MP (monitor coprocessor) and TS (task switched) bits of CR0 are set. |
| 8 | Double fault | Abort | SIGSEGV | When the CPU detects an exception while trying to invoke the handler for a prior exception, the two exceptions can be handled serially. If, however, the processor cannot handle them serially, it signals this exception instead. |
| 9 | Coprocessor segment overrun | Abort | SIGFPE | Raised in protected mode if the 80386 detects a page or segment violation while transferring the middle portion of a coprocessor operand to the NPX. |
| 10 | Invalid TSS | Fault | SIGSEGV | Occurs if during a task switch the new TSS is invalid. |
| 11 | Segment not present | Fault | SIGBUS | A reference was made to a segment not present in memory. |
| 12 | Stack segment | Fault | SIGBUS | The instruction attempted to exceed the stack segment limit, or the segment identified by ss is not present in mem. |

| | | | | |
|---|---|---|---|---|
| 13 | General protection | Fault | SIGSEGV | All protection violations that do not cause anohter exception cause this exception. |
| 14 | Page fault | Fault | SIGSEGV | This exception occurs when CPU finds the page is not present in memory. |
| 15 | Reserved | | | |
| 16 | Coprocessor error | Fault | SIGFPE | The CPU reports this exception when it detects a signal from the 80387 on the 80386's ERROR# input pin. |

# 5.5  traps.c

## 5.5.1  Functions

traps.c program mainly includes the C functions invoked by the low level code in asm.s file. These C functions is usually used to show the error location and error code of the exception. In the current Linux system, these C functions is usually called the bottom halfs, because the quantity of their code or processing speed is relately slow, thus, they may be interrupted by other code in the kernel.

The die() function is used to display detailed error information about the exception on the console. The trap_init() function at the end of traps.c program, used to initialize interrupt vectors for the CPU or hardware exceptionsis and enable the interrupt requests, is invoked in the init/main.c program. When read this program, it's better to refer to asm.s program at the same time.

## 5.5.2  Comments

Program 5-3 linux/kernel/traps.c

```
1  /*
2   *  linux/kernel/traps.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  /*
8   * 'Traps.c' handles hardware traps and faults after we have saved some
9   * state in 'asm.s'. Currently mostly a debugging-aid, will be extended
10  * to mainly kill the offending process (probably by giving it a signal,
11  * but possibly by killing it outright if necessary).
12  */
13  #include <string.h>        // Declares one type and several functions for manipulating strings.
14
15  #include <linux/head.h>    // define a simple structure of segment descript and a few constans.
16  #include <linux/sched.h>
17  #include <linux/kernel.h>
18  #include <asm/system.h>
19  #include <asm/segment.h>
```

```
20 #include <asm/io.h>
21
   // The following lines define three inline assembly macros. A compound statement in parentheses
   // may appear inside an expression in GNU C. This allows you to declare variables within an
   // expression. The last variable (__res) is the output value of the macro. Refers to the
   // descriptions for inline assembly instructions.
   //
   // Get one byte from address 'addr' in segment 'seg'.
   // Arguments: seg - segment selector; addr - offset in the segment.
   // %0 - eax (__res); %1 - eax (seg); %2 - offset (*(addr)).
22 #define get_seg_byte(seg,addr) ({ \
23 register char __res; \
24 __asm__("push %%fs;mov %%ax,%%fs;movb %%fs:%2,%%al;pop %%fs" \
25         :"=a" (__res):"" (seg), "m" (*(addr))); \
26 __res;})
27
   // Get a long word (4-byte) from address 'addr' in segment 'seg'.
   // Arguments: seg - segment selector; addr - offset in the segment.
   // %0 - eax (__res); %1 - eax (seg); %2 - offset (*(addr)).
28 #define get_seg_long(seg,addr) ({ \
29 register unsigned long __res; \
30 __asm__("push %%fs;mov %%ax,%%fs;movl %%fs:%2,%%eax;pop %%fs" \
31         :"=a" (__res):"" (seg), "m" (*(addr))); \
32 __res;})
33
   // Get the value of fs register (selector).
   // %0 - eax (__res).
34 #define _fs() ({ \
35 register unsigned short __res; \
36 __asm__("mov %%fs,%%ax":"=a" (__res):); \
37 __res;})
38
39 int do_exit(long code);                         // At L102 in kernel/exit.c
40
41 void page_exception(void);                      // page_fault. At L14 in mm/page.s
42
   // The following lines define the prototypes of interrupt handling procedure. They were
   // used to setup the corresponding interrupt gate descriptors in function trap_init().
43 void divide_error(void);                        // int0 (kernel/asm.s,L19)。
44 void debug(void);                               // int1 (kernel/asm.s,L53)。
45 void nmi(void);                                 // int2 (kernel/asm.s,L57)。
46 void int3(void);                                // int3 (kernel/asm.s,L61)。
47 void overflow(void);                            // int4 (kernel/asm.s,L65)。
48 void bounds(void);                              // int5 (kernel/asm.s,L69)。
49 void invalid_op(void);                          // int6 (kernel/asm.s,L73)。
50 void device_not_available(void);                // int7 (kernel/system_call.s,L148)。
51 void double_fault(void);                        // int8 (kernel/asm.s,L97)。
52 void coprocessor_segment_overrun(void);         // int9 (kernel/asm.s,L77)。
53 void invalid_TSS(void);                         // int10 (kernel/asm.s,L131)。
54 void segment_not_present(void);                 // int11 (kernel/asm.s,L135)。
55 void stack_segment(void);                       // int12 (kernel/asm.s,L139)。
56 void general_protection(void);                  // int13 (kernel/asm.s,L143)。
57 void page_fault(void);                          // int14 (mm/page.s,L14)。
```

```
58 void coprocessor_error(void);                    // int16 (kernel/system_call.s,L131)。
59 void reserved(void);                             // int15 (kernel/asm.s,L81)。
60 void parallel_interrupt(void);                   // int39 (kernel/system_call.s,L280)。
61 void irq13(void);                                // int45 (kernel/asm.s,L85)。
62
   // This function is used to print the interrupt name, error code, and the messages related to
   // the calling process including eip, eflags, esp, base address and limits of segments,
   // process pid, task number and 10 bytes opcodes. If the stack is in user data segment, then
   // the 16 bytes of stack is printed too.
63 static void die(char * str,long esp_ptr,long nr)
64 {
65         long * esp = (long *) esp_ptr;
66         int i;
67
68         printk("%s: %04x\n\r",str,nr&0xffff);
69         printk("EIP:\t%04x:%p\nEFLAGS:\t%p\nESP:\t%04x:%p\n",
70                 esp[1],esp[0],esp[2],esp[4],esp[3]);
71         printk("fs: %04x\n",_fs());
72         printk("base: %p, limit: %p\n",get_base(current->ldt[1]),get_limit(0x17));
73         if (esp[4] == 0x17) {
74                 printk("Stack: ");
75                 for (i=0;i<4;i++)
76                         printk("%p ",get_seg_long(0x17,i+(long *)esp[3]));
77                 printk("\n");
78         }
79         str(i);                          // Get task number of the current process. (L159, sched.h)
80         printk("Pid: %d, process nr: %d\n\r",current->pid,0xffff & i);
81         for(i=0;i<10;i++)
82                 printk("%02x ",0xff & get_seg_byte(esp[1],(i+(char *)esp[0])));
83         printk("\n\r");
84         do_exit(11);                     /* play segment exception */
85 }
86
   // The following functions are invoked by relating interrupt handlers.
87 void do_double_fault(long esp, long error_code)
88 {
89         die("double fault",esp,error_code);
90 }
91
92 void do_general_protection(long esp, long error_code)
93 {
94         die("general protection",esp,error_code);
95 }
96
97 void do_divide_error(long esp, long error_code)
98 {
99         die("divide error",esp,error_code);
100 }
101
102 void do_int3(long * esp, long error_code,
103                 long fs,long es,long ds,
104                 long ebp,long esi,long edi,
105                 long edx,long ecx,long ebx,long eax)
```

```
106 {
107         int tr;
108
109         __asm__("str %%ax":"=a"(tr):""(0)); // Get task register's value into variable tr.
110         printk("eax\t\tebx\t\tecx\t\tedx\n\r%8x\t%8x\t%8x\t%8x\n\r",
111               eax, ebx, ecx, edx);
112         printk("esi\t\tedi\t\tebp\t\tesp\n\r%8x\t%8x\t%8x\t%8x\n\r",
113               esi, edi, ebp, (long) esp);
114         printk("\n\rds\tes\tfs\ttr\n\r%4x\t%4x\t%4x\t%4x\n\r",
115               ds, es, fs, tr);
116         printk("EIP: %8x    CS: %4x   EFLAGS: %8x\n\r", esp[0], esp[1], esp[2]);
117 }
118
119 void do_nmi(long esp, long error_code)
120 {
121         die("nmi", esp, error_code);
122 }
123
124 void do_debug(long esp, long error_code)
125 {
126         die("debug", esp, error_code);
127 }
128
129 void do_overflow(long esp, long error_code)
130 {
131         die("overflow", esp, error_code);
132 }
133
134 void do_bounds(long esp, long error_code)
135 {
136         die("bounds", esp, error_code);
137 }
138
139 void do_invalid_op(long esp, long error_code)
140 {
141         die("invalid operand", esp, error_code);
142 }
143
144 void do_device_not_available(long esp, long error_code)
145 {
146         die("device not available", esp, error_code);
147 }
148
149 void do_coprocessor_segment_overrun(long esp, long error_code)
150 {
151         die("coprocessor segment overrun", esp, error_code);
152 }
153
154 void do_invalid_TSS(long esp, long error_code)
155 {
156         die("invalid TSS", esp, error_code);
157 }
158
```

```
159 void do_segment_not_present(long esp,long error_code)
160 {
161         die("segment not present",esp,error_code);
162 }
163
164 void do_stack_segment(long esp,long error_code)
165 {
166         die("stack segment",esp,error_code);
167 }
168
169 void do_coprocessor_error(long esp, long error_code)
170 {
171         if (last_task_used_math != current)
172                 return;
173         die("coprocessor error",esp,error_code);
174 }
175
176 void do_reserved(long esp, long error_code)
177 {
178         die("reserved (15,17-47) error",esp,error_code);
179 }
180
    // This is the initialization function of interrupt exception used to setup various interrupt
    // gates.
    // The main difference between set_trap_gate() and set_system_gate() is that the former's
    // privilege level is set to 0, and the later's is set to 3. Thus the handler of breakpoint
    // int3, overflow and bounds interrupt can be invoked by any programs. This two macros were
    // implemented in include/asm/system.h file.
181 void trap_init(void)
182 {
183         int i;
184
185         set_trap_gate(0,&divide_error);
186         set_trap_gate(1,&debug);
187         set_trap_gate(2,&nmi);
188         set_system_gate(3,&int3);        /* int3-5 can be called from all */
189         set_system_gate(4,&overflow);
190         set_system_gate(5,&bounds);
191         set_trap_gate(6,&invalid_op);
192         set_trap_gate(7,&device_not_available);
193         set_trap_gate(8,&double_fault);
194         set_trap_gate(9,&coprocessor_segment_overrun);
195         set_trap_gate(10,&invalid_TSS);
196         set_trap_gate(11,&segment_not_present);
197         set_trap_gate(12,&stack_segment);
198         set_trap_gate(13,&general_protection);
199         set_trap_gate(14,&page_fault);
200         set_trap_gate(15,&reserved);
201         set_trap_gate(16,&coprocessor_error);
    // The following loop first batch installing the vectors of int 17-48 to be the handler of
    // 'reserved'. Each hardware initial routine will reinstall its handler descriptor later.
202         for (i=17;i<48;i++)
203                 set_trap_gate(i,&reserved);
```

```
     // Installing the trap gates of copprocessor and parallel device. Enabling the interrupt
     // request of PIC chips.
204         set_trap_gate(45,&irq13);              // Trap gate of coprocessor.
205         outb_p(inb_p(0x21)&0xfb,0x21);         // Enable IRQ2 of 8259A master chip.
206         outb(inb_p(0xA1)&0xdf,0xA1);           // Enable IRQ13 of slave chip.
207         set_trap_gate(39,&parallel_interrupt); // Trap gate of parallel device.
208 }
209
```

## 5.5.3 Information

### 5.5.3.1 GNU gcc Inline assembly

This is the first time we encounter the inline assembly instructions in C language. Because we have very few chance to use it in normal application programming, thus, here, we present a simple description about its syntax, enough to get the ability to understand the kernel source code. Refers to GNU gcc manual or an article "Using Inline Assembly with gcc" for a detailed description of it.

The following is the basic format of the inline assembly with input and output arguments:

```
asm("statements"
    : output_registers
    : input_registers
    : clobbered_registers);
```

"statements" is the location you write assembly instructions, 'output_registers' (or called output operands) indicates where the output values were in the registers after finished executing the assembly instructions. Here, these registers are associated with C variables or an address in memory seperately. 'input_registers' (or called input operands) represents the input values each indicated register has before start running assembly instructions. They may associated with C expressions or constant values. 'clobbered_registers' is a set of registers used to tell gcc can no longer count on the values it loaded into these registers. When we use inline assembly, the last three lines of the format can be omited if not used.

Let's explain the using method with an example. The following lines are a piece of code selected begin from line 22 from previous program (Program 5-3). We have rearranged it to have a better layout.

```
01  #define get_seg_byte(seg,addr) \
02  ({ \
03  register char __res; \
04  __asm__("push %%fs; \
05          mov %%ax,%%fs; \
06          movb %%fs:%2,%%al; \
07          pop %%fs" \
08          :"=a" (__res) \
09          :"" (seg),"m" (*(addr))); \
```

<u>10</u>  __res;})

---

This piece of code defines a macro function using the inline assembly. Usually the most convenient way to use inline assembly is to encapsulate them in macros. A compound statement in parentheses (sentences in bracket parentheses) may appear inside an expression in GNU C. This allows us to declare variables within an expression. The value of the variable __res at the last line is the output value of the expression.

The first line (L01) defines macro's identifier with parameters givien in the parentheses. L02 is the beginning of a compound statement. L03 declares a register variable '__res' in the compound statement. The '__asm__' at L04 starts the inline assembly statements. Instructions within L04-07 are AT&T format assembly sentences.

L08 is the output operand section, contains one operand here. The '__res' is the C expression of the output operand. Output operand expressions must be lvalues and write-only. '"a"' is the operand constraint, saying that a eax register is required. The '=' in '=a' indicates that the operand is an output. All output operands' constraints must use '='. Each operand is described by an operand-constraint string followed by the C expression in paretheses. A colon separates the assembler statements from the first output operand, and another separates the last output operand from the first input operand, if any. Commas are used to separate output operands and separate input operands.

If there are no output operands, and there are input operands, then there must be two consecutive colons surrounding the place where the output operands would go.

L09 is the input operand section, contains two operands here. The 'seg' is the C expression for the first input operand while the '*(addr)' is that of the second. The '""' is an ellipsis of '"0"' for input operand 1, means that it must occupy the same location as operand 0. A digit in constraint is allowed only in an input operand, and it must refer to an output operand. In order to refer to the operands in the inline assembler statements, all operands are ordered begin from zero (%0) with the first operand in output operands or input operands if output operands are not existed. For example the input operand '*(addr)' is refered to with '%2' in the assembler statements. The '"m"' is the constraint for the second input operand, means that the second operand is a memory effective address.

The input operands must be read-only. When the assembler instruction has a read-write operand, you must logcally split its function into two separate operands, one input operand and one write-only output operand.

Now, let's take a look for the meaning of assembler statements at L04-07. The first instruction push the segment register fs onto stack; The second instruction, then, store the eax into fs, that is store the value of 'seg' into fs; The third instruction save the value at 'fs:(*(addr))' into register al. When finished executing the assembler instructions, this value is stored into variable '__res' in according to the output operand constraint. So the whole macro means: 'seg' indicates a segment in memory, 'addr' is an offset address in the segment. The macro returns one byte value pointed by 'seg' and 'addr' in memory.

Notice that in the inline assembly, you must prefix registers with '%%' instead of just '%'. Because as GCC parses along for '%0' and '%1' and so on, it would interpret '%edx' as a '%e' parameter, and that is non-existent, and cause an error.

Let's see another example below.

```
01  asm("cld\n\t"
02      "rep\n\t"
03      "stol"
04      : /* no outputs */
05      : "c"(count-1), "a"(fill_value), "D"(dest)
06      : "%ecx", "%edi");
```

The '\ns' and '\t's are there so the '.s' file that gcc generates and hands to gas comes out right when you have got multiple statements per asm. It's really meant for issuing instructions for which there is no equivalent in C and don't touch the registers.

L01-03 are the normal assembler instructions used to clear the direction of string operating and repeat saving value. L04 indicates that is inline assembly has no output operand. L05 is the input operands: the value of 'count-1' will be stored in register ecx (operand constraint is "c"), eax with 'fill_value' and edi with 'dest'. Why make GCC do it instead of doing it ourselves? Because GCC, in its register allocating, might be able to arrange for, say, 'fill_value' to already be in eax. If this is in a loop, it might be able to preserve eax through the loop, and save a movl once per loop. L06 specifies GCC that it can no longer count on the values it loaded into ecx or edi to be valid. This doesn't mean they will be reloaded for certain. This is the clobberlist.

Seem strange? Well, it really helps when GCC do optimizing, when GCC can know exactly what you are doing with the registers before and after. It folds your assembly code into the code it generates and then optimizes. Table 5-2 is a list of operand constraints that we often use.

Table 5-2 Operand constraints for inline assembly

| Code | Description | Code | Description |
|------|-------------|------|-------------|
| a | eax | m | memory effective address |
| b | ebx | o | memory address with offset |
| c | ecx | I | constant value (0 to 31). |
| d | edx | J | constant value (0 to 63). |
| S | esi | K | constant value (0 to 255). |
| D | edi | L | constant value (0 to 65535). |
| q | Dynamicaly addressed register. (eax, ebx, ecx or edx) | M | constant value (0 to 3). |
| r | Dynamicaly addressed register. (eax, ebx, ecx, edx, esi, edi) | N | one byte constant value (0 to 255). |
| g | general register (eax, ebx, ecx, edx or variable in memory). | O | constant value (0 to 31). |
| A | eax and edx combined into a 64-bit integer (use long longs) | n | immediate value known at compile time. |

Its possible to let GCC choose the registers for you as in the following example:

```
01  asm("leal (%1, %1, 4), %0"
02      : "=r"(y)
03      : "0"(x));
```

The first assembly statement (L01) "leal (r1, r2, 4), r3" means r1 + r2*4 --> r3. This example multiplies 'x' by 5 really quickly. Now, we could have specified, say eax. But unless we really need a specific register (like when using rep movsl or rep stosl, which are hardcoded to use ecx, edi, and esi), why not let GCC pick an available one? So when GCC generates the output code for GAS, %0 will be replaced by the register it picked.

Notice that if you want to prevent an asm instruction from beging deleted, moved or combined by GCC, you can write keyword 'volatile' after the 'asm'. For example:

asm volatile (...);

The following example is a little longer one selected from file include/string.h.

```
//// Compare first count bytes between string str1 and str2.
// Arguments: cs - Str1，ct - Str2，count - Compared byte count.
// %0 - eax(__res) ret value; %1 - edi(cs) point to str1; %2 - esi(ct) point to str2; %3 - ecx(count)。
// Return 1 if str1 > str2; 0 if str1=str2; -1 if str1 < str2.
extern inline int strncmp(const char * cs,const char * ct,int count)
{
register int __res ;                // register variable.
__asm__("cld\n"                     // Clear direction bit in eflags.
        "1:\tdecl %3\n\t"           // count--。
        "js 2f\n\t"                 // if count < 0 then jump forward to label 2.
        "lodsb\n\t"                 // Get str2's char ds:[esi] into al and esi++
        "scasb\n\t"                 // Compare al and str1's char es:[edi] and edi++
        "jne 3f\n\t"                // Forward jump to label 3 if not equal.
        "testb %%al,%%al\n\t"       // a NULL char ?
        "jne 1b\n"                  // Backward jump to label 1 if not a NULL char.
        "2:\txorl %%eax,%%eax\n\t"  // Clear eax and return if tt's a NULL char.
        "jmp 4f\n"                  // Forward jump to label 4 and return.
        "3:\tmovl $1,%%eax\n\t"     // eax = 1
        "jl 4f\n\t"                 // Return 1 if str2 < str1,
        "negl %%eax\n"              // else return -1.
        "4:"
        :"=a" (__res):"D" (cs),"S" (ct),"c" (count)
        :"si","di","cx");
return __res;                       // return result.
}
```

# 5.6 system_call.s

## 5.6.1 Function

In the Linux system, user uses interrupt int 0x80 and the function number in register eax to invoke the function provided by kernel. These operating system internal functions are called the system call functions. Usually, the application program does not invoke the system call interrupt directly. It uses the relating interfacing function provided by general library such as libc.a. For example, when creating a new process, it usually uses function fork() in the library. The function fork() in library will implement the procedure of doing int 0x80 and return the result to user.

For all system calls, the kernel lists the function pointers in the table sys_call_table[] (in include/linux/sys.h) ordered by the function number. In the procedure of int 0x80, it invokes the system function based on the function number provied by the user in eax.

The system_call.s program is mainly used to implement the interface of system calls with user application program and detect and process the signals of the calling process. In addition, two system functions (sys_execve for executing a program and sys_fork for forking a new process) and several interrupt handling procedures are also implemented in this program. The interrupt handling procedures includes coprocessor error (int 16), device not found (int 7), timer interrupt (int 32), hard disk interrupt (int 46) and floppy device interrupt (int 38).

For the software interrupts (system calls, coprocessor_error, device_not_available), their handling procedures are basically identical: push some arguments for the C functions to call, invoke the C function, when return, check the signal map of the current process and process the siganl with the lowest value and reset the siganl processed. The C functions for all system calls are scattered in the kernel source tree. They are assembled to be a function pointer table in header file include/linux/sys.h.

For the hardware interrupt requests, the handling procedure first sends 'end of interrupt' (EOI) signal to the 8259A chip and then invokes the corresponding C functions. The handling of timer interrupt also needs to process the signal map of current process.

The handling procedure of system call (int 0x80) in this program can be considered as a interface to the real processing procedures in each corresponding C function (the bottom halfs). At the beginning of this program, it first checks the validity of the system call number in register eax and then saves all the registers used. The kernel uses the segment register ds and es for kernel data segment, fs for user data segment as default. Afterwards, it invokes the relating C functions according to the address table of C functions and saves the return value onto the stack when returns from the C funtion.

After that, it checks the state of the calling process. If the operation of the relating C function has changed the state of the calling process or the time for the calling process was used up (counter ==0), then it invokes the schedule function to switch to the other process. Although it uses 'jmp _schedule' instruction to execute the schedule function, the CPU control will be passed to here finally after the kernel switches back to the current process. Because

the program pushed the return address (ret_from_sys_call) onto the stack before jump to the schedule function.

The code begin from label 'ret_from_sys_call' is used to do some ending operation for this interrupt execution. It checks the current process's type to decide if to end the interrupt handling immediately. If the process is a kernel process (e.g. process 0 or 1), it ends the handling right away. Otherwise it will handle the sigals of the process further more. If the process does have siganls based on its signal map, it will handling the signal by calling function do_signal().

At the end, the system_call.s program restores the registers saved and exits from this interrupt handling procedure. If the process has signal，then the control will pass to the signal handling function first, and after that, back to the calling program.

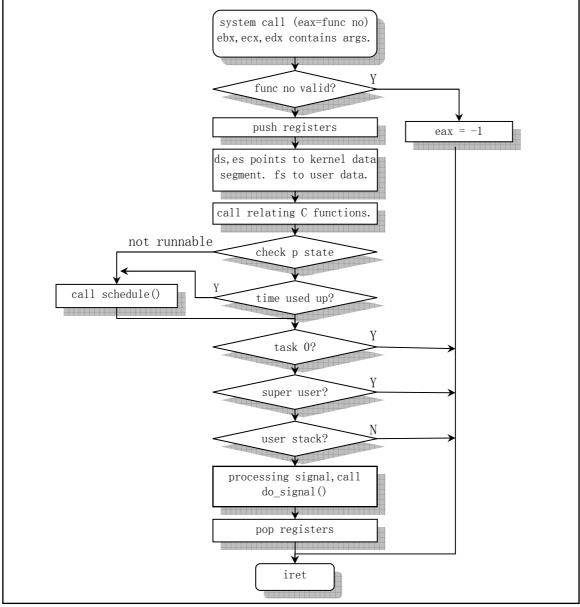The flow chart of the system call interrupt processing is illustrated in Figure 5-5.



Figure 5-5 The flow chart of system call interrupt.

## 5.6.2 Comments

Program 5-4 linux/kernel/system_call.s

```
1  /*
2   *  linux/kernel/system_call.s
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  /*
8   *  system_call.s  contains the system-call low-level handling routines.
9   * This also contains the timer-interrupt handler, as some of the code is
10  * the same. The hd- and flopppy-interrupts are also here.
11  *
12  * NOTE: This code handles signal-recognition, which happens every time
13  * after a timer-interrupt and after each system call. Ordinary interrupts
14  * don't handle signal-recognition, as that would clutter them up totally
15  * unnecessarily.
16  *
17  * Stack layout in 'ret_from_system_call':
18  *
19  *       0(%esp) - %eax
20  *       4(%esp) - %ebx
21  *       8(%esp) - %ecx
22  *       C(%esp) - %edx
23  *      10(%esp) - %fs
24  *      14(%esp) - %es
25  *      18(%esp) - %ds
26  *      1C(%esp) - %eip
27  *      20(%esp) - %cs
28  *      24(%esp) - %eflags
29  *      28(%esp) - %oldesp
30  *      2C(%esp) - %oldss
31  */
32
33  SIG_CHLD       = 17          # Number of SIG_CHLD signal. (child stopped or end)
34
35  EAX            = 0x00        # The relative offset of registers on stack.
36  EBX            = 0x04
37  ECX            = 0x08
38  EDX            = 0x0C
39  FS             = 0x10
40  ES             = 0x14
41  DS             = 0x18
42  EIP            = 0x1C
43  CS             = 0x20
44  EFLAGS         = 0x24
45  OLDESP         = 0x28        # When privilege changes.
46  OLDSS          = 0x2C
47
    # The relative offset of some fields in task_struct. Refer to include/linux/sched.h.
48  state  = 0              # these are offsets into the task-struct.
```

```
49 counter = 4              # ticks of task running (Time slice), decreasing.
50 priority = 8             # counter=priority when start running.
51 signal  = 12             # signal map, each bit associated with one kind of signal.
52 sigaction = 16           # MUST be 16 (=len of sigaction). The offset of sigaction struct.
53 blocked = (33*16)        # signal blocked.
54
55 # offsets within sigaction        # Refer to include/signal.h
56 sa_handler = 0                    # signal handler.
57 sa_mask = 4                       # siganl mask.
58 sa_flags = 8                      # signal flags.
59 sa_restorer = 12                  # Refer to kernel/signal.c.
60
61 nr_system_calls = 72   # Total system call number in Linux 0.11
62
63 /*
64  * Ok, I get parallel printer interrupts while using the floppy for some
65  * strange reason. Urgel. Now I just ignore them.
66  */
67 .globl _system_call,_sys_fork,_timer_interrupt,_sys_execve
68 .globl _hd_interrupt,_floppy_interrupt,_parallel_interrupt
69 .globl _device_not_available, _coprocessor_error
70
   # Wrong system call number (out of range).
71 .align 2                          # align to 4 bytes bound.
72 bad_sys_call:
73        movl $-1,%eax              # set eax = -1
74        iret
   # Do scheduling. Refer to L104 in kernel/sched.c
75 .align 2
76 reschedule:
77        pushl $ret_from_sys_call   # Push the return address (L101).
78        jmp _schedule
   #### int 0x80 -- Entry point for Linux system call (int 0x80, eax = function number).
79 .align 2
80 _system_call:
81        cmpl $nr_system_calls-1,%eax   # Check the function number in eax.
82        ja bad_sys_call
83        push %ds                   # Save registers.
84        push %es
85        push %fs
86        pushl %edx                 # ebx,ecx,edx contains arguments for C function if any.
87        pushl %ecx                 # push %ebx,%ecx,%edx as parameters
88        pushl %ebx                 # to the system call
89        movl $0x10,%edx            # set up ds,es to kernel space
90        mov %dx,%ds
91        mov %dx,%es
92        movl $0x17,%edx            # fs points to local data space
93        mov %dx,%fs
   # The calling address = _sys_call_table + %eax * 4
   # _sys_call_table is an address arry of system call functions.
94        call _sys_call_table(,%eax,4)
95        pushl %eax                 # Push the return value of function.
96        movl _current,%eax         # Get the address of current process --> eax.
```

```
     # The following lines check the state of current process. If it isn't in running state
     # ( state !=0) or it is in running state but its time is used up, then do rescheduling.
 97          cmpl $0,state(%eax)        # state
 98          jne reschedule
 99          cmpl $0,counter(%eax)      # counter
100          je reschedule
     # The following code handling the siganl of current process.
101 ret_from_sys_call:
     # Check to see whether the current process is task 0, if true then return immediately.
     # The '_task' on L103 is the array name of task[] in C. Direct reference to an array name
     # equals refering to task[0].
102          movl _current,%eax         # task[0] cannot have signals
103          cmpl _task,%eax
104          je 3f                      # Jump forward to label 3.
     # Checking the segment selector of the calling process to see whether the process is kernel
     # code. If true then return immediately. If false, then handle the signal of the process.
     # The following code compares the old code segment selector with user code segment selector
     # 0x0f (RPL = 3, LDT, segment 1 (Local code segment)).
105          cmpw $0x0f,CS(%esp)        # was old code segment supervisor ?
106          jne 3f
     # If the stack selector is not 0x17 (that is it's not a user stack), then return.
107          cmpw $0x17,OLDSS(%esp)     # was stack segment = 0x17 ?
108          jne 3f
     # The following code (L109-120) is used to handle process signal. First it gets the signal
     # map from the task_struct (32-bit, one bit for a signal). Then blocks the siganls with the
     # the signal mask. After that it gets the signal with lowest value and reset the siganl in
     # signal map. And, finally, invokes the do_signal() with the siganl as its one argument.
     # do_siganl() is implemented at L82 in kernel/signal.c and has 13 arguments.
109          movl signal(%eax),%ebx     # Get signal map --> ebx.
110          movl blocked(%eax),%ecx    # Get signal block mask --> ecx.
111          notl %ecx                  # Inverse each bit.
112          andl %ebx,%ecx             # Obtain the permissive signal map.
113          bsfl %ecx,%ecx             # Scan the map from bit 0. ecx will contains seq. number.
114          je 3f                      # If no signal then return.
115          btrl %ecx,%ebx             # Reset the signal (ebx contains the orignal map).
116          movl %ebx,signal(%eax)     # Save the siganl map --> current->signal.
117          incl %ecx                  # Adjust the signal value start from 1 (1-32).
118          pushl %ecx                 # Push the siganl value as one of the arguments for
119          call _do_signal            # function do_signal().
120          popl %eax                  # Pop out the siganl value.
121 3:       popl %eax
122          popl %ebx
123          popl %ecx
124          popl %edx
125          pop %fs
126          pop %es
127          pop %ds
128          iret
129
     #### int 16 -- Handling the error of coprocessor.
     # Jump to invoke C function math_error() at L82 in kernel/math/math_emulate.c. When return,
     # it will continue executing start from ret_from_sys_call.
130 .align 2
```

```
131 _coprocessor_error:
132         push %ds
133         push %es
134         push %fs
135         pushl %edx
136         pushl %ecx
137         pushl %ebx
138         pushl %eax
139         movl $0x10,%eax          # ds,es points to kernel space.
140         mov %ax,%ds
141         mov %ax,%es
142         movl $0x17,%eax          # fs points to user space.
143         mov %ax,%fs
144         pushl $ret_from_sys_call # Push the return address.
145         jmp _math_error          # Invoking the C function math_error().
146
    #### int7 -- Device not available (Coprocessor not available).
    # The CPU interprets the pattern 11011B in the first five bits of an instruction as an opcode
    # intended for a coprocessor. Instructions thus marked are called ESCAPE or ESC instructions.
    # If the CR0's EM (Emulation Mode) flag is set, it means that we should emulate the coprocessor
    # functions. If the CPU finds EM set when executing an ESC instruction, it signals this
    # exception, giving the exception handler an opportunity to emulate the ESC instruction.
    # The TS (Task switched) flag of CR0 is set while CPU does task switching. So we can depend
    # on this flag to know whether the contents of the copprocess matches with the current task
    # and do some updates of the state of the coprocessor, if necessary. Refer to L77 in sched.c.
147 .align 2
148 _device_not_available:
149         push %ds
150         push %es
151         push %fs
152         pushl %edx
153         pushl %ecx
154         pushl %ebx
155         pushl %eax
156         movl $0x10,%eax          # ds,es points to kernel data space.
157         mov %ax,%ds
158         mov %ax,%es
159         movl $0x17,%eax          # fs points to user space.
160         mov %ax,%fs
161         pushl $ret_from_sys_call # Push the return address.
162         clts                     # clear TS so that we can use math
163         movl %cr0,%eax
164         testl $0x4,%eax          # EM (math emulation bit)
    # If this exception is not caused by EM flag, then update the state of coprocessor for new task.
165         je _math_state_restore   # Invoke the function at L77 in kernel/sched.c.
166         pushl %ebp
167         pushl %esi
168         pushl %edi
    # Else we should emulate the coprocessor in program kernel/math/math_emulate.c.
169         call _math_emulate
170         popl %edi
171         popl %esi
172         popl %ebp
```

```
173          ret                        # ret(jump) to label ret_from_sys_call.
174
     #### int32 -- (int 0x20) Timer interrupt handler.
     # The Linux system timer is set to the frequency of 100Hz (at L5 in include/linux/sched.h)
     # by initializing the 8253 timer chip at L406 in kernel/sched.c. Thus the system time unit
     # jiffies is 10 milliseconds.
     # This exception handler first increases jiffies by 1, sends the EOI command to nofiy
     # the end of the interrupt for interrupt controller (PIC), then invokes the C function
     # do_timer() with current privilege level (CPL) as its argument.
175 .align 2
176 _timer_interrupt:
177          push %ds                    # save ds,es and put kernel data space
178          push %es                    # into them. %fs is used by _system_call
179          push %fs
180          pushl %edx                  # we save %eax,%ecx,%edx as gcc doesn't
181          pushl %ecx                  # save those across function calls. %ebx
182          pushl %ebx                  # is saved as we use that in ret_sys_call
183          pushl %eax
184          movl $0x10,%eax             # ds,es points to kernel space.
185          mov %ax,%ds
186          mov %ax,%es
187          movl $0x17,%eax             # fs points to user space.
188          mov %ax,%fs
189          incl _jiffies
     # As the kernel does not initialize PIC to use auto EOI, so here we send EOI command manually.
190          movb $0x20,%al              # EOI to interrupt controller #1
191          outb %al,$0x20              # OCW2 send to port 0x20 of PIC
     # The following code is used to get CPL from old cs selector and pushed onto stack as a
     # argument for C function do_timer().
192          movl CS(%esp),%eax
193          andl $3,%eax                # %eax is CPL (0 or 3, 0=supervisor)
194          pushl %eax
     # The jobs of do_timer() is to do task switching, time conunting (at L305 in sched.c).
195          call _do_timer              # 'do_timer(long CPL)' does everything from
196          addl $4,%esp                # task switching to accounting ...
197          jmp ret_from_sys_call
198
     #### This is the entry point of system call sys_execve. It invokes the C function
     # do_execve() (L182 of fs/exec.c) with the old eip as its argument.
199 .align 2
200 _sys_execve:
201          lea EIP(%esp),%eax
202          pushl %eax                  # Push old eip as an argument of do_execve().
203          call _do_execve
204          addl $4,%esp                # Get rid of the argument.
205          ret
206
     #### Entry point for system call sys_fork() used to fork a new process.
     # It first call the C function find_empty_process() to get a new pid number, then invokes
     # the C function copy_process() to duplicate the contents of the process for child.
207 .align 2
208 _sys_fork:
209          call _find_empty_process   # L135 in kernel/fork.c
```

```
210         testl %eax,%eax
211         js 1f
212         push %gs                    # These are part of the arguments for copy_process.
213         pushl %esi
214         pushl %edi
215         pushl %ebp
216         pushl %eax
217         call _copy_process          # L68 in kernel/fork.c
218         addl $20,%esp               # Get rid of the arguments pushed.
219 1:      ret
220
```

    #### int 46 -- (int 0x2E) Hard disk interrupt handler responding hardwrae int request IRQ14.
    # This interrupt is signaled when the request hard disk operating finished or encounters an
    # error during processing. Refer to kernel/blk_drv/hd.c file.
    # This handler first send EOI to the slave 8259A chip, then gets and checks the C function
    # pointer in variable do_hd. If the pointer is null, it invokes unexpected_hd_interrupt() to
    # show the error message, else it will invokes the C function using the pointer in do_hd.
    # Register edx is used to store the pointer of C functions, which is one of read_intr(),
    # write_intr() or unexpected_hd_interrupt().

```
221 _hd_interrupt:
222         pushl %eax
223         pushl %ecx
224         pushl %edx
225         push %ds
226         push %es
227         push %fs
228         movl $0x10,%eax             # ds,es points to kernel space.
229         mov %ax,%ds
230         mov %ax,%es
231         movl $0x17,%eax             # fs points to user space.
232         mov %ax,%fs
```
    # As the kernel does not initialize PIC to use auto EOI, so here we send EOI command manually.
```
233         movb $0x20,%al
234         outb %al,$0xA0             # EOI to interrupt controller #1 (slave 8259A chip)
235         jmp 1f                     # give port chance to breathe
236 1:      jmp 1f
```
    # do_hd is defined as a C func pointer. When the pointer in do_hd is stored in edx, it will
    # be set to null.
```
237 1:      xorl %edx,%edx
238         xchgl _do_hd,%edx
```
    # Check to see if do_hd is null orignally. If null, then edx points to unexpected_hd_interrupt()
    # implemented at L237 in file kernel/blk_drv/hdc.
```
239         testl %edx,%edx
240         jne 1f
241         movl $_unexpected_hd_interrupt,%edx
242 1:      outb %al,$0x20            # Send EOI to master 8259A chip.
243         call *%edx                # "interesting" way of handling intr.
244         pop %fs                   # The above calls the C function pointed by edx.
245         pop %es
246         pop %ds
247         popl %edx
248         popl %ecx
249         popl %eax
```

```
250         iret
251
    #### int 38 -- (int 0x26) Interrupt handler for floppy drive. Responding int request IRQ6.
    # The processing way is similar as for hard disk exception int 46. Refer to blk_drv/floppy.c
    # This handler first send EOI to the slave 8259A chip, then gets and checks the C function
    # pointer in variable do_floppy. If the pointer is null, it invokes nexpected_floppy_interrupt()
    # to show the error message, else it will invokes the C function using the pointer in do_floppy.
    # Register eax is used to store the pointer of C functions, which is one of rw_interrupt(),
    # seek_interrupt(), recal_interrupt(), reset_interrupt() or unexpected_floppy_interrupt().
252 _floppy_interrupt:
253         pushl %eax
254         pushl %ecx
255         pushl %edx
256         push %ds
257         push %es
258         push %fs
259         movl $0x10,%eax          # ds,es points to kernel space.
260         mov %ax,%ds
261         mov %ax,%es
262         movl $0x17,%eax          # fs points to user space.
263         mov %ax,%fs
    # As the kernel does not initialize PIC to use auto EOI, so here we send EOI command manually.
264         movb $0x20,%al
265         outb %al,$0x20           # EOI to interrupt controller #1 (master 8259A chip)
    # do_floppy is defined as a C func pointer. When the pointer in do_hd is stored in edx, it
    # will be set to null.
266         xorl %eax,%eax
267         xchgl _do_floppy,%eax
    # Check to see if do_floppy is null orignally. If null, then sets register edx to point to
    # unexpected_floppy_interrupt(),which implemented in file kernel/blk_drv/floppy.c.
268         testl %eax,%eax
269         jne 1f
270         movl $_unexpected_floppy_interrupt,%eax
271 1:      call *%eax                # "interesting" way of handling intr.
272         pop %fs
273         pop %es
274         pop %ds
275         popl %edx
276         popl %ecx
277         popl %eax
278         iret
279
    #### int 39 -- (int 0x27) Interrupt handler for parallel port. Responding request IRQ7.
    # Not yet implemented in this kernel.
280 _parallel_interrupt:
281         pushl %eax
282         movb $0x20,%al
283         outb %al,$0x20
284         popl %eax
285         iret
```

## 5.6.3  Infomation

### 5.6.3.1 The representation of addressing methods in GNU assembly

The GNU as or inline assembly in gcc uses AT&T syntax. The formal 32 bits addressing format is as follows:

    AT&T: immed32(basepointer, indexpointer, indexscale)
    Intel: [basepointer + indexpointer*indexscal + immed32]

The calculation method of this format is : immed32 + basepointer + indexpointer * indexscale. In an application we needn't write out all these fields. But one of the immed32 and basepointer must exists. The following are some examples.

◆ Addressing a particular C variable:

AT&T: _booga                         Intel: [_booga]

Notice that the underscore ('_') is how you get at static C variables from assembler. This only works with global variables. Otherwise, you can use extended asm to have variables preloaded into registers for you.

◆ Addressing what a register points to:

AT&T: (%eax)                         Intel: [eax]

◆ Addressing a variable offset by a value in a registr:

AT&T: _variable(%eax)                Intel: [eax + _variable]

◆ Addressing a value in an array of integers (scaling up by 4):

AT&T: _array(,%eax,4)                Intel: [eax*4 + _array]

◆ You can also do offsets with the immediate value:

C code: *(p+1) where p is a char *

AT&T:   1(%eax) where eax has the value of p        Intel: [eax + 1]

◆ You can do some simple math on the immediate value:

AT&T: _struct_pointer + 8            Intel: [_struct_pointer + 8]

◆ Addressing a particular char in an array of 8-character records. eax holds the number of the record desired. ebx has the wanted char's offset within the record.

AT&T: _array(%ebx,%eax,8)            Intel: [ebx + eax*8 + _array]

### 5.6.3.2  Method for adding new system call

To add a new system call to the Linux 0.11 kernel, we need to do the following things.

First we should write a C function processing the system call, for example, sys_sethostname(). This function is used to change the hostname of the system. Usually, we can put this function in file kernel/sys.c. You can certainly put your system call function in a proper location in the kernel source tree. The following program list is one kind of implementation of sys_sethostname(), selected from Linux kernel 0.12, as an example here.

Notice that this function uses variable thisname of structure utsname type in it. This variable is defined in function sys_uname() internally. So we need to move it out from sys_uname() (at L218-220 in sys.c) to be as a global difinition. The hostname field of the is set as an constant value "linux .0" in the file.

```
#define MAXHOSTNAMELEN 8
int sys_sethostname(char *name, int len)
{
        int     i;

        if (!suser())
                return -EPERM;
        if (len > MAXHOSTNAMELEN)
                return -EINVAL;
        for (i=0; i < len; i++) {
                if ((thisname.nodename[i] = get_fs_byte(name+i)) == 0)
                        break;
        }
        if (thisname.nodename[i]) {
                thisname.nodename[i>MAXHOSTNAMELEN ? MAXHOSTNAMELEN : i] = 0;
        }
        return 0;
}
```

Then, we need to allocate a number to this new system call in the header file include/unistd.h. Since Linux 0.11 has 72 system calls (number from 0-71), we append this new one at the end of it. So, we can add a macro difinition for it as follows (properly inserted after line 131). We also need to add the new system call function prototype in the same header file .

```
// In include/unistd.h file.
// New system call number. properly inserted after line 131.
#define __NR_sethostname     72

// The function prototype. properly inserted just after line 251.
int sethostname(char *name, int len);
```

After that, we need to modify the function address table sys_call_table[] of the system calls to append the new function address. This array type table is located in file include/linux/sys.h. The modification is shown as following. Notice that the system call function pointers must be arranged in sequence in the array according to their function number.

```
// In file include/linux/sys.h

extern int sys_sethostname();

// An array table of function pointers for system calls. The new function name is appended at
// the end of the array.
fn_ptr sys_call_table[] = { sys_setup, sys_exit,sys_fork, sys_read,
 ...,
sys_setreuid, sys_setregid, sys_sethostname };
```

Now, we should tell the kernel that the total number of system calls has increased. This can be achived by increasing the number of nr_system_calls by one at L61 in file system_call.s.

At the end, we can recompile to build a kernel Image file which contains our new system call function. To test this new function in our application programs, we may add a interface program in our C library as follows. Refer to chapter 12 for a description of a similar program.

```
#define __LIBRARY__
#include <unistd.h>

_syscall2(int, sethostname, char *, name, int, len);
```

# 5.7  mktime.c

## 5.7.1  Function

This program continas one function mktime(), which is used only in init/main.c file for kernel initialization. It calculates the time passed in seconds beginning from the first day in 1970 to the moment of kernel power on.

## 5.7.2  Comments

Program 5-5 linux/kernel/mktime.c

```
 1  /*
 2   *  linux/kernel/mktime.c
 3   *
 4   *  (C) 1991  Linus Torvalds
 5   */
 6
 7  #include <time.h>
 8
 9  /*
10   * This isn't the library routine, it is only used in the kernel.
11   * as such, we don't care about years<1970 etc, but assume everything
12   * is ok. Similarly, TZ etc is happily ignored. We just do everything
13   * as easily as possible. Let's find something public for the library
14   * routines (although I think minix times is public).
15   */
16  /*
17   * PS. I hate whoever though up the year 1970 - couldn't they have gotten
18   * a leap-year instead? I also hate Gregorius, pope or no. I'm grumpy.
19   */
20  #define MINUTE 60                      // Seconds in a minute.
21  #define HOUR (60*MINUTE)               // Seconds in an hour.
22  #define DAY (24*HOUR)                  // Seconds in a day.
23  #define YEAR (365*DAY)                 // Seconds in a year.
24
25  /* interestingly, we assume leap-years */
    // Bellow is an array of time begin from the first day of a year to the start of each month.
26  static int month[12] = {
```

```
27          0,
28          DAY*(31),
29          DAY*(31+29),
30          DAY*(31+29+31),
31          DAY*(31+29+31+30),
32          DAY*(31+29+31+30+31),
33          DAY*(31+29+31+30+31+30),
34          DAY*(31+29+31+30+31+30+31),
35          DAY*(31+29+31+30+31+30+31+31),
36          DAY*(31+29+31+30+31+30+31+31+30),
37          DAY*(31+29+31+30+31+30+31+31+30+31),
38          DAY*(31+29+31+30+31+30+31+31+30+31+30)
39 };
40
   // Calculate the time between January 1 1970 and system power on in seconds.
   // The argument tm is a structure, each field of it is filled with values collected from CMOS.
41 long kernel_mktime(struct tm * tm)
42 {
43          long res;
44          int year;
45
46          year = tm->tm_year - 70;           // Years from 1970 to now. suck 2k problem.
47 /* magic offsets (y+1) needed to get leapyears right.*/
   // Seconds for these years plus one more day for each leapyear and month in current year.
48          res = YEAR*year + DAY*((year+1)/4);
49          res += month[tm->tm_mon];
50 /* and (y+2) here. If it wasn't a leap-year, we have to adjust */
51          if (tm->tm_mon>1 && ((year+2)%4))
52               res -= DAY;
53          res += DAY*(tm->tm_mday-1);        // Plus seconds of days in current month.
54          res += HOUR*tm->tm_hour;           // Plus seconds of hours in the current day.
55          res += MINUTE*tm->tm_min;          // Plus seconds of minutes in the current hour.
56          res += tm->tm_sec;                 // Plus seconds in the current minute.
57          return res;                        // Finally get the seconds passed from 1970.1.1
58 }
59
```

### 5.7.3  Information

#### 5.7.3.1 Leapear calculation method

The year is a leapyear if it can be divided exactly by 4 and cannot be divided exactly by 100, or can be divided exactly by 400.

# 5.8  sched.c

## 5.8.1  Function

Program sched.c is mainly concerned with task switching management in the kernel. The primary functions include sleep_on(), wakeup(), schedule() and do_timer(). In addition, it also contains several C functions for system calls and other "trivial" routines. To facilitate the

kernel programming, Linus put the timing routines about floppy drives here too.

schedule() function first checks all tasks to activate every task that has received signal(s). For each task in the kernel, it checks alarm time value in the task's structure. If the alarm time has expired (alarm<jiffies), it set the SIGALRM signal in the signal map of the task and then reset the alarm value to zero. After ignoring the signals blocked, if the task still has other signal and its running state is in TASK_INTERRUPTIBLE, then schedule() will change its state to TASK_RUNNING.

After that, the schedule() function enters into task switching part. It selects the task to be run based on the time slice and the priority scheduling mechanism. First, it checks the counter value of all tasks in ready state. The counter value represents the remaining running time before the task be switched. It selects the task that has the largest value among all tasks and switches to it by using the switch_to() function. If the counter value of all ready to run tasks are zero, which means their running time were used out, then the scheduler will reset the counter value of all tasks, including the tasks at other states, based on their priority value. Finally, it repeat the processing above to selcet one task to run.

Other two functions need to be mentioned are sleep_on() and wake_up(). sleep_on() is used by a process to enter into sleep state by itself. wake_up() is used to change the state of a task to TASK_RUNNING manually. Although small in size, these two functions are somehow hard to understand. Here we discuss them with the help of graph illustrating.

The main function of sleep_on() is to put the current process into sleep state (TASK_UNINTERRUPTIBLE) in a waiting queue for a period of time when the resources it requests is busy or not in memory. A task in TASK_UNINTERRUPTIBLE state need to be wake up manually by using the function wake_up(). The way the task be put into a waiting queue is to use the temporary pointer 'tmp' in the function to be as the relationship links between each waiting tasks in the queue.

The sleep_on() function manipulates three task pointers: '*p', 'tmp' and 'current'. '*p' is a argument of the function indicating the header pointer of the waiting queue. It represents, for eaxmple, the i_wait pointer in the file system waiting queue, the buffer_wait pointer in buffer operating waiting queue. 'tmp' is used as a link pointer in a waiting queue. 'current' the the pointer of current task. We can illustrated the variation of these pointers in memory in Figure 5-6. The long line in the figure represents the bytes sequence in memory.
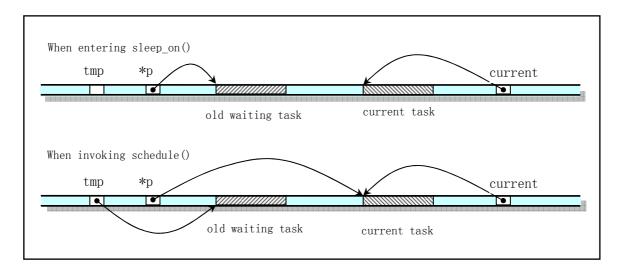
Figure 5-6 The variation of pointer in memory

When just entering the sleep_on() function, the queue header pointer '*p' points to the task already waiting on the queue, or is null if the queue is empty at the beginning. By manipulating these pointers, before invoking the scheduler, the queue header pointer '*p' pointed to the current process while the 'tmp' pointed to the task already waiting in the queue. In other words, the new task (current) is inserted into the queue formed by the links of 'tmp'. Thus, by using the temporary pointer 'tmp' in the function for each waiting task in the queue, a implicit waiting queue is formed. From the Figure 5-7, we can better understand the constructing process of the queue. The figure shows the case when inserting the third task into the waiting queue.
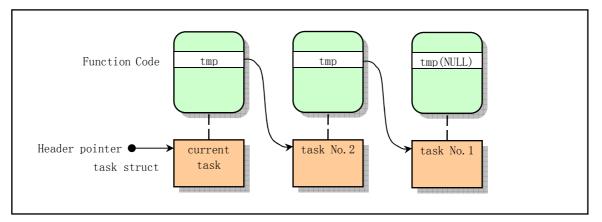


Figure 5-7 The implicit waiting queue for sleep_on()

After the current process was inserted into the queue, sleep_on() invokes the scheduler to switch to run another task. When the task pointed by queue header pointer is waken up, it continues to execute the code bellow 'schedule()' in sleep_on() and wakes up all the other tasks waited in the queue.

The wake_up() function is used to wakes up the tasks waiting for resources on a waiting queue. It is a general waking up function. In some situations, for example to read the data block on a disk, since any task on the waiting queue can be waked up first, the pointer to the waked up task should be set to null. Thus, when the task lately entered into the queue is waked up and continues to run the code in the sleep_on(), it needn't to wake the task with null pointer.

The other function, interruptible_sleep_on(), has the similar structure of sleep_on(). The only difference between them is that the former set the state of the sleeping task into TASK_INTERRUPTIBLE and checks whether there are any other tasks on the queue. In the Linux kenrel 0.12, These two functions are merged into one subroutine and uses the indicated task state as its argument to differentiate these two cases.

When reading the following code, its better to refer to the comments in include/kernel/sched.h file.

## 5.8.2 Comments

Program 5-6 linux/kernel/sched.c

1 /*

```
2   *  linux/kernel/sched.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  /*
8   * 'sched.c' is the main kernel file. It contains scheduling primitives
9   * (sleep_on, wakeup, schedule etc) as well as a number of simple system
10  * call functions (type getpid(), which just extracts a field from
11  * current-task
12  */
13 #include <linux/sched.h>
14 #include <linux/kernel.h>
15 #include <linux/sys.h>
16 #include <linux/fdreg.h>
17 #include <asm/system.h>
18 #include <asm/io.h>
19 #include <asm/segment.h>
20
21 #include <signal.h>
22
   // Get the bit mask (binary value) corresponding to signal nr in signal bitmap. The signal
   // value range is 1 to 32. e.g.: nr=5, then the bit mask = 1<<(5-1) = 16 = 00010000b.
23 #define _S(nr) (1<<((nr)-1))
   // All signals are blockable except signal SIGKILL and SIGSTOP.
24 #define _BLOCKABLE (~(_S(SIGKILL) | _S(SIGSTOP)))
25
   // Show the task number, pid, state and the kernel stack free space of the nr-th task in
   // the task array.
26 void show_task(int nr,struct task_struct * p)
27 {
28         int i,j = 4096-sizeof(struct task_struct);
29
30         printk("%d: pid=%d, state=%d, ",nr,p->pid,p->state);
   // The task's kernel stack is located after the task's struct and on the same page. The loop
   // bellow detects the zero value bytes after the task's struct on the page.
31         i=0;
32         while (i<j && !((char *)(p+1))[i])
33                 i++;
34         printk("%d (of %d) chars free in kernel stack\n\r",i,j);
35 }
36
   // Show the task number, pid, state and the kernel stack free space of all tasks.
37 void show_stat(void)
38 {
39         int i;
40
   // NR_TASKS (64) is the total task number the system can hold at one time. It's also the
   // size of array task[] and is defined at L4 in file include/kernel/sched.h.
41         for (i=0;i<NR_TASKS;i++)
42                 if (task[i])
43                         show_task(i,task[i]);
44 }
```

```
45
      // The input frequency of timer 8253 chip is 1.193180MHz. Linux wishes the timer interrupt
      // is 100 Hz, that is generating an interrupt request signal IRQ0 every 10ms (every tick).
      // Thus the LATCH is the initial value for the channel 0 of the 8253 chip. Refer to L407.
46 #define LATCH (1193180/HZ)
47
48 extern void mem_use(void);         // [??] Not used.
49
50 extern int timer_interrupt(void);  // L176 in file kernel/system_call.s
51 extern int system_call(void);      // L80 in kernel/system_call.s
52
      // Define a task union. The task descriptor and the task kernel stack is on the same
      // memory page. Each task uses its own kernel stack when run in kernel mode. The top of
      // the stack is at the end of the page.
53 union task_union {
54         struct task_struct task;
55         char stack[PAGE_SIZE];
56 };
57
      // Defined the default data of the first task (task 0 or called the idle process).
58 static union task_union init_task = {INIT_TASK,};
59
      // Variable 'jiffies' is the time ticks begin from the power on of the system. The type-
      // qualifier 'volatile' is used to suppress optimization that could otherwise occur to the
      // volatile object (jiffies) by the compiler. Defined at L139 in include/linux/sched.h.
60 long volatile jiffies=0;
61 long startup_time=0;                                  // Power on time (secs) from 1970.1.1
62 struct task_struct *current = &(init_task.task); // Pointer points to current task.
63 struct task_struct *last_task_used_math = NULL;
64
65 struct task_struct * task[NR_TASKS] = {&(init_task.task), };
66
      // User stack space. It will be used as a user mode stack for task 0 (the idle process)
      // after the kernel finished initialization. 4 Kb in size.
67 long user_stack [ PAGE_SIZE>>2 ] ;
68
      // Define the pointer of user_stack ss:esp. Refer to L23 in boot/head.s. esp will point
      // the last item of user_stack[].
69 struct {
70         long * a;
71         short b;
72         } stack_start = { & user_stack [PAGE_SIZE>>2] , 0x10 };
73 /*
74  * 'math_state_restore()' saves the current math information in the
75  * old math state array, and gets the new ones from the current task
76  */
      // Save the math state array for the switched out task and restore the new ones for the
      // new switched in task.
77 void math_state_restore()
78 {
79         if (last_task_used_math == current) // Return if task does not change.
80                 return;
      // Save the math state if the last task used the math coprocessor. Before issue any other
```

```
     // math instructions, we need to execute 'fwait' first.
81           __asm__("fwait");
82           if (last_task_used_math) {
83                   __asm__("fnsave %0"::"m" (last_task_used_math->tss.i387));
84           }
     // Memorize the current task which is to use the math coprocessor. If the current task
     // used math before then restore its math state, else it means this is the first time
     // the current task is going to use math, then issue initializing command to coprocessor
     // and set the 'used_math' flag in the task descriptor.
85           last_task_used_math=current;
86           if (current->used_math) {            // If used before, then resotore its state.
87                   __asm__("frstor %0"::"m" (current->tss.i387));
88           } else {                             // Else the firt time to use.
89                   __asm__("fninit"::);         // Issue the initilizing command.
90                   current->used_math=1;        // Set use math flag.
91           }
92 }
93
94 /*
95  *  'schedule()' is the scheduler function. This is GOOD CODE! There
96  * probably won't be any reason to change this, as it should work well
97  * in all circumstances (ie gives IO-bound processes good response etc).
98  * The one thing you might take a look at is the signal-handler code here.
99  *
100 *   NOTE!!  Task 0 is the 'idle' task, which gets called when no other
101 * tasks can run. It can not be killed, and it cannot sleep. The 'state'
102 * information in task[0] is never used.
103 */
104 void schedule(void)
105 {
106         int i,next,c;
107         struct task_struct ** p;         // A pointer to the pointer of task descriptor.
108
109 /* check alarm, wake up any interruptible tasks that have got a signal */
110
     // Check 'alarm' begin from the last task in the task pointer array.
111         for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
112                 if (*p) {
     // If the task set the alarm time before and now is expired, then issue the SIGALRM signal
     // to the task and clear the value of 'alarm'. SIGLARM's default opeararion is to terminate
     // the task.
113                         if ((*p)->alarm && (*p)->alarm < jiffies) {
114                                 (*p)->signal |= (1<<(SIGALRM-1));
115                                 (*p)->alarm = 0;
116                         }
     // If this task has got signal(s) that can not be blocked (SIGKILL and SIGSTOP) or not blocked,
     // and the task is in state TASK_INTERRUPTIBLE, then change the state of the task to ready.
     // Code '~(_BLOCKABLE & (*p)->blocked)' is used to ignore signals blocked. But the signal
     // SIGKILL and SIGSTOP cannot be blocked or ignored.
117                         if (((*p)->signal & ~(_BLOCKABLE & (*p)->blocked)) &&
118                         (*p)->state==TASK_INTERRUPTIBLE)
119                                 (*p)->state=TASK_RUNNING;
120                 }
```

121
122 */* this is the scheduler proper: */
123
124          while (1) {
125                  c = -1;
126                  next = 0;
    // The following code is also doing the loop begin from the last item of the task pointer
    // array. It ignores the null pointers and compares the running time counter value of each
    // task in ready state. The counter value is the remain time slice a task need to run before
    // being switched. The pointer 'next' will point to the task which has the largest counter
    // value and 'c' will contain the counter value.
127                  i = NR_TASKS;
128                  p = &task[NR_TASKS];
129                  while (--i) {
130                          if (!*--p)
131                                  continue;
132                          if ((*p)->state == TASK_RUNNING && (*p)->counter > c)
133                                  c = (*p)->counter, next = i;
134                  }
    // If there is one task its counter value is great than zero, then the program exits the
    // loop starting from L124 and switches to this task at L141.
135                  if (c) break;
    // Otherwise updates each task's counter value based on their priority and repeat loop.
    // The counter value is set to equal to 'counter /2 + priority' with each own values .
136                  for(p = &LAST_TASK ; p > &FIRST_TASK ; --p)
137                          if (*p)
138                                  (*p)->counter = ((*p)->counter >> 1) +
139                                                  (*p)->priority;
140          }
    // Now we switch to task 'next' selected in the loop above. Since the 'next' is initialized
    // to 0, if there is no other tasks could be switched to, the scheduler will run the task 0
    // (idle process) at system idle time. The task 0 invokes system call pause(), which in turn
    // invokes schedule() again.
141          switch_to(next);
142 }
143
    // C function for system call pause(). Change the state of current task to be in waiting
    // state: TASK_INTERRUPTIBLE, and rescheduling.
    // The pause function will cause the invoking process to sleep until a signal is received
    // that either terminates it or causes it to call a signal-catching fuction. This function
    // only returns when a signal was caught and the signal-catching function returned. In this
    // case puase() returns -1, and errno is set to EINTR. So the full functions of pause() are
    // not yet wholly implemented in kernel 0.11 (until version 0.95).
144 int sys_pause(void)
145 {
146          current->state = TASK_INTERRUPTIBLE;
147          schedule();
148          return 0;
149 }
150
    // Due to resources lacking or waiting for other events, the current process initiatively
    // changes itself into sleep state TASK_UNINTERRUPTIBLE and be putted onto a implicit waiting
    // queue by invoking this function. Processes at this state can be back into running state

```
     // only through the explicit wake up operation using function wake_up() (at L188).
     // The argument '*p' is the header pointer of a waiting queue.
151  void sleep_on(struct task_struct **p)
152  {
153          struct task_struct *tmp;
154
     // The object pointed by a pointer can be NULL, but a valid pointer itself can not be null.
155          if (!p)
156                  return;
     // if the current process is init task (task 0), then panic (impossible!).
157          if (current == &(init_task.task))
158                  panic("task[0] trying to sleep");
     // Set the temporary pointer 'tmp' points to the task already waited on the queue (if any) and
     // set the header pointer of queue points to the current task. Variable 'tmp' is created on
     // the task's kernel stack and used as implicit links for the waiting queue.
159          tmp = *p;
160          *p = current;
     // Change the state of current task to sleep and do rescheduling.
161          current->state = TASK_UNINTERRUPTIBLE;
162          schedule();
     // Only when the task is changed back to running state again by explicitly invoking the
     // wake_up() function, will the following code be continued. At this moment, the task is
     // explicitly waked up. While the tasks on the queue are waiting for the same resources or
     // events to occur, they will be waked up together by the code bellow.
163          if (tmp)                              // Wake up the previous task on the queue.
164                  tmp->state=0;
165  }
166
     // Set the current task into interruptible sleep state, and insert into a waiting queue pointed
     // by pointer '*p'. The processing way of this function is similar as function sleep_on().
167  void interruptible_sleep_on(struct task_struct **p)
168  {
169          struct task_struct *tmp;
170
171          if (!p)
172                  return;
173          if (current == &(init_task.task))
174                  panic("task[0] trying to sleep");
175          tmp=*p;
176          *p=current;
     // Change the state of current task to sleep and do rescheduling.
177  repeat: current->state = TASK_INTERRUPTIBLE;
178          schedule();
     // If the queue has other waiting tasks and the header pointer of the queue doest not point
     // to the current task, then change the state of the task pointed by the header pointer to
     // running state and let the current task sleep again.
     // When the '*p' does not point to the current task, it means there is at least one new task
     // be inserted into the waiting queue after the current task. Thus we should wake up these
     // tasks too. This also indicates that the last task inserted in the queue will be the first
     // one to use the resources waited for.
179          if (*p && *p != current) {
180                  (**p).state=0;
181                  goto repeat;
```

```
182             }
```
    // The following line should be '*p = tmp', that is let the header pointer to point to the
    // remain tasks on the queue. Otherwise, the tasks previously inserted on the queue are
    // erroneously ignored. In addition, the code on L192 should be deleted at the same time.
```
183             *p=NULL;
184             if (tmp)
185                     tmp->state=0;
186 }
187
```
    // Waking up the task pointed by '**p'.
```
188 void wake_up(struct task_struct **p)
189 {
190             if (p && *p) {
191                     (**p).state=0;    // Change the task to running state.
192                     *p=NULL;
193             }
194 }
195
196 /*
197  * OK, here are some floppy things that shouldn't be in the kernel
198  * proper. They are here because the floppy needs a timer, and this
199  * was the easiest way of doing it.
200  */
```
    // The following code (L201-262) deals with the timing operation of floppy drives. It is
    // recommended to read this block of code with program floppy.c in the next chapter simultaneously
    // or after. The time unit tick equals 1/100 second.

    // The array contains the task descriptor pointer waited for the floppy motor speed up.
    // Array item 0 - 3 corresponding to the floppy drive 0 - 3 (A - D).
```
201 static struct task_struct * wait_motor[4] = {NULL,NULL,NULL,NULL};
```
    // Before data can be transferred to or from the diskette, the drive motor must be brought
    // up to speed. The bellow array will filled with the motor speed up delay time (in ticks)
    // for each floppy drives. For most floppy drives, the spin-up time is 50 ticks (0.5 second).
```
202 static int  mon_timer[4]={0,0,0,0};
```
    // The bellow array hold the remainning spin time for each floppy drive before stopped.
    // The default time is 10000 ticks (100 seconds).
```
203 static int moff_timer[4]={0,0,0,0};
```
    // Digital Output Register (DOR) contains the drive select and motor enable bits, a reset bit
    // and a DMAGATE# bit. This variable holds the current value of DOR sent to de floppy controller.
    // Bit 7-4: Each bit controls the motor for drive D-A seperately. 1 - start; 0 - stop.
    // Bit 3  : 1 - Enable DMA and interrupt request; 0 - Disable DMA and interrupt request.
    // Bit 2  : 1 - Enable floppy controller; 0 - Disable floppy controller.
    // Bit 1-0: 00 -11 selects drive A to D.
```
204 unsigned char current_DOR = 0x0C;       // Enable DMA and int, enable controller.
205
```
    // Decide the spin-up (motor on) delay time.
    // The argument nr (0 - 3) indicates the diskette drive.
    // Returns the spin-up delay time in ticks.
```
206 int ticks_to_floppy_on(unsigned int nr)
207 {
208             extern unsigned char selected;   // Drive selected flag. (L122,blk_drv/floppy.c)
209             unsigned char mask = 0x10 << nr; // Corresponds to the motor bits in DOR.
210                                              // The high 4-bit of mask is the motor on flag.
```

```
211         if (nr>3)
212             panic("floppy_on: nr>3");  // Contains no more than 4 floppy drives.
    // First we preset the motor turn off delay time (100 seconds). Then merge the contents
    // of current_DOR into variable mask.
213         moff_timer[nr]=10000;              /* 100 s = very big :-) */
214         cli();                             /* use floppy_off to turn it off */
215         mask |= current_DOR;
    // If there is no selected drive at this moment, then we select the indicated one (nr).
216         if (!selected) {
217             mask &= 0xFC;
218             mask |= nr;
219         }
    // If the current value of DOR is different with that of the requested, then we output the
    // new value in 'mask' to the floppy controller. Here, if the drive motor has not been turned
    // on, then we should set the drive's motor spin-up time to 0.5 s (HZ/2 = 50 ticks), else
    // we reset it to 2 ticks to satisfy the requests for decreasing first in function
    // do_floppy_timer(). Finally, we save the current DOR into variable 'current_DOR'.
220         if (mask != current_DOR) {
221             outb(mask,FD_DOR);
222             if ((mask ^ current_DOR) & 0xf0)
223                 mon_timer[nr] = HZ/2;
224             else if (mon_timer[nr] < 2)
225                 mon_timer[nr] = 2;
226             current_DOR = mask;
227         }
228         sti();
229         return mon_timer[nr];             // Return the motor spin-up delay time.
230 }
231
    // Waiting for the motor to spin-up,then return.
    // Argument 'nr' specifies the floppy drive.
    // Set the motor spin-up time if the DOR is different with that of the request. Sleep to
    // wait for the motor to speed up. In the handle of timing interrupt (L331 in do_timer()),
    // the spin-up time (mon_timer[nr]) will be decreased and checked, until the time is up.
232 void floppy_on(unsigned int nr)
233 {
234         cli();                            // disable interrupt.
    // While the motor spin-up time is not over, we put the current task into motor waiting
    // queue with uninterruptible state.
235         while (ticks_to_floppy_on(nr))
236             sleep_on(nr+wait_motor);
237         sti();                            // Enable interrupt.
238 }
239
    // Set the motor spin-down (turn off) delay time for the specified floppy drive.
    // If we do not explicitly turn off the motor within 3 seconds using this function, the
    // motor will be turned off after 100 seconds.
240 void floppy_off(unsigned int nr)
241 {
242         moff_timer[nr]=3*HZ;
243 }
244
    // Timer function for floppy drives.
```

```
      // Updating the value of motor spin-up and turn off delay time in mon_timer[] and moff_timer[]
      // seperately. This function will be invoked in the C function do_timer() of system timing
      // interrupt. Thus the kernel invokes the function for each ticks to update the delay time.
      // If the delay time for turning off the motor is up, the corresponding motor flag bit in
      // DOR will be cleared.
245 void do_floppy_timer(void)
246 {
247         int i;
248         unsigned char mask = 0x10;
249
250         for (i=0 ; i<4 ; i++,mask <<= 1) {
251                 if (!(mask & current_DOR))            // Skip if not the specified motor.
252                         continue;
253                 if (mon_timer[i]) {
      // If the motor spin-up time is over then wakes up the corresponding waiting task(s).
254                         if (!--mon_timer[i])
255                                 wake_up(i+wait_motor);
256                 } else if (!moff_timer[i]) {
      // If the motor spin-down delay time is over, then reset the corresponding motor flag in
      // DOR, save and update the DOR in floppy controller.
257                         current_DOR &= ~mask;
258                         outb(current_DOR,FD_DOR);
259                 } else
260                         moff_timer[i]--;            // Decreasing the turn-off delay time.
261         }
262 }
263
264 #define TIME_REQUESTS 64            // Maximum timers.
265
      // 定时器链表结构和定时器数组。
266 static struct timer_list {
267         long jiffies;                     // Time to delay (ticks).
268         void (*fn)();                     // Handler, the address of a function to be called.
269         struct timer_list * next;     //
270 } timer_list[TIME_REQUESTS], * next_timer = NULL; // The header pointer of the timer_list.
271
      // Add timer function.
      // The floppy driver (floppy.c) uses this function to control the state of its motor.
      // Arguments: jiffies - ticks to delay; *fn() - the handler when timer expired.
272 void add_timer(long jiffies, void (*fn)(void))
273 {
274         struct timer_list * p;
275
      // The timer handler can not be null.
276         if (!fn)
277                 return;
278         cli();
      // If the timer value jiffies <=0, then we should process the timer handler right away and
      // needn't insert it into the timer list.
279         if (jiffies <= 0)
280                 (fn)();
281         else {
      // Otherwise, we will add it to the timer list. First search to get a empty item in the
```

```
        // array timer_list[].If there is no empty item availabe, then the kernel goes panic.
282               for (p = timer_list ; p < timer_list + TIME_REQUESTS ; p++)
283                       if (!p->fn)
284                               break;
285               if (p >= timer_list + TIME_REQUESTS)
286                       panic("No more time requests free");
        // Insert the request timer into and to be the first item of the timer list. The timer
        // header pointer points to it.
287               p->fn = fn;
288               p->jiffies = jiffies;
289               p->next = next_timer;
290               next_timer = p;
        // Now we sort the items of the timer list, the order depends on the delay value, from
        // short to long dely time. During the sorting, the time value of a item need to minus the
        // sum time of all previous items. Thus, we need only to check the first item in the timer
        // list to determine if there is a timer expired.
291               while (p->next && p->next->jiffies < p->jiffies) {
292                       p->jiffies -= p->next->jiffies;
293                       fn = p->fn;
294                       p->fn = p->next->fn;
295                       p->next->fn = fn;
296                       jiffies = p->jiffies;
297                       p->jiffies = p->next->jiffies;
298                       p->next->jiffies = jiffies;
299                       p = p->next;
300               }
301         }
302       sti();
303 }
304
    // The C function of the handler of system timer interrupt invoked at L176 in system_call.s
    // Argument 'cpl' is the current privilege level. If it is zero, means CPU is executing the
    // kernel code. This function updates the timer and counting the time used by the current
    // process. When the process uses out its time slice and is running its user's code, then
    // the CPU will do switching
305 void do_timer(long cpl)
306 {
307       extern int beepcount;      // beep count of the speaker. (L697 in chr_drv/console.c)
308       extern void sysbeepstop(void);  // Stop the speaker. (L691 in chr_drv/console.c)
309
    // If the beep count is reached, then stop the speaker. The way is to send command to port 0x61,
    // reset bit 0 and bit 1. Bit 0 controls the channel 2 operation of the 8253 timer chip.
    // Bit 1 controls the speaker.
310       if (beepcount)
311               if (!--beepcount)
312                       sysbeepstop();
313
    // Updating the resource usage. If the argument 'cpl' is 0, means the CPU is executing kernel
    // code, then increasing the kernel counting time stime of the process, else increasing the
    // utime of the process. [By reading the kernel code, we can see that Linus named the kernel
    // code as supervisor's code. Refer to Linus' comments at L193 in file system_call.s ]
314       if (cpl)
315               current->utime++;
```

```
316            else
317                    current->stime++;
318
```

        // By checking the timer list header pointer, we can decide if there is timer exist in the
        // timer list. If there is one, then minus the time value by 1 to the first timer. If the
        // value equals zero, means the timer is expired, then we invoke its handler and finally
        // remove the item from the timer list.

```
319            if (next_timer) {                    // timer list header (see L270).
320                    next_timer->jiffies--;
321                    while (next_timer && next_timer->jiffies <= 0) {
322                            void (*fn)(void);    // A definition of a func pointer ! :(
323
324                            fn = next_timer->fn;
325                            next_timer->fn = NULL;
326                            next_timer = next_timer->next;
327                            (fn)();              // invoking the timer handler.
328                    }
329            }
```

        // If the digital output register (DOR) of the floppy controller contains the motor turning
        // on bit(s), then invokes the floppy timer handler 'do_floppy_timer()' (L245).

```
330            if (current_DOR & 0xf0)
331                    do_floppy_timer();
```

        // If the time slice of the current process is not yet expired, then return. Otherwise, if
        // the CPU is executing kernel code before entered into this interrupt handler, we also
        // return right away. This means that the Linux kernel code cannot be preempted. Otherwise
        // we invoke the scheduler to switch the current process.

```
332            if ((--current->counter)>0) return;
333            current->counter=0;
334            if (!cpl) return;
335            schedule();
336 }
337
```

        // System call function - Set a timer that will expire at a specified time in the future.
        // When the timer expires, the SIGALRM signal is generted for the process. If we ignore or
        // don't catch this signal, its default action is to terminate the process.
        // Argument 'seconds' is the number of seconds in the future when the signal should be generated.
        // If it greater than 0, then it is set in the alarm field of the process descriptor, else
        // a zero is returned.

```
338 int sys_alarm(long seconds)
339 {
340            int old = current->alarm;
341
342            if (old)
343                    old = (old - jiffies) / HZ;
344            current->alarm = (seconds>0)?(jiffies+HZ*seconds):0;
345            return (old);
346 }
347
```

        // Get process identification, returns  the process ID of the current process - pid.

```
348 int sys_getpid(void)
349 {
350            return current->pid;
351 }
```

352
```
    // Returns the process ID of the parent of the current process - ppid.
353 int sys_getppid(void)
354 {
355         return current->father;
356 }
357
    // Get user identity, returns the real user ID of the current process - uid.
358 int sys_getuid(void)
359 {
360         return current->uid;
361 }
362
    // Returns the effective user ID of the current process - euid.
363 int sys_geteuid(void)
364 {
365         return current->euid;
366 }
367
    // Get group identity, returns the real group ID of the current process - gid.
368 int sys_getgid(void)
369 {
370         return current->gid;
371 }
372
    // Returns the effective group ID of the current process - egid.
373 int sys_getegid(void)
374 {
375         return current->egid;
376 }
377
    // Change process priority, adds 'increment' to the nice value for the current process.
    // A large nice value means a low priority. A negative value will increase the privority.
378 int sys_nice(long increment)
379 {
380         if (current->priority-increment>0)
381                 current->priority -= increment;
382         return 0;
383 }
384
    // Schedule initialization routine.
385 void sched_init(void)
386 {
387         int i;
388         struct desc_struct * p;   // descriptor structure, refer to include/linux/head.h
389
390         if (sizeof(struct sigaction) != 16)    // sigaction is a signal handling struct.
391                 panic("Struct sigaction MUST be 16 bytes");
    // Setup the init task's TSS and LDT descriptors in the global descriptor table (GDT).
    // Refer to lines begin from 65 in file include/asm/system.h
392         set_tss_desc(gdt+FIRST_TSS_ENTRY,&(init_task.task.tss));
393         set_ldt_desc(gdt+FIRST_LDT_ENTRY,&(init_task.task.ldt));
    // Clear the task pointer array and the descriptors in GDT. Notice that the action begin
```

```
        // from the second item in the task array, so the init task pointer is still there.
394         p = gdt+2+FIRST_TSS_ENTRY;
395         for(i=1;i<NR_TASKS;i++) {
396             task[i] = NULL;
397             p->a=p->b=0;
398             p++;
399             p->a=p->b=0;
400             p++;
401         }
402 /* Clear NT, so that we won't have troubles with that later on */
        // The NT (Nested Task) flag in the system EFLAGS controls the chaining of interrupted and
        // called tasks. It influences the operation of the IRET instruction. If NT is set, the IRET
        // instruction will cause the task switch. Here reset the NT flag.
403         __asm__("pushfl ; andl $0xffffbfff,(%esp) ; popfl");
        // Load the init task's TSS selector into task register TR from the GDT. Load the init
        // task's LDT selector into local descriptor table register LDTR from the GDT. Notice that
        // the LDTR is only explicitly loaded once at here, all the subsequent loading of LDTR will
        // be performed by CPU automatically according to the LDT field in TSS.
404         ltr(0);
405         lldt(0);
        // The following code is used to initialize the channnel 0 of the 8253 timer chip. Select
        // mode 3, binary counting method. The output pin of channel 0 is connected to the IRQ0 of
        // of the master PIC 8259A chip. It issues a IRQ0 signal every 10 ms. The constant LATCH
        // is the initial counting value set to the channel 0, refer to L46.
406         outb_p(0x36,0x43);              /* binary, mode 3, LSB/MSB, ch 0 */
407         outb_p(LATCH & 0xff , 0x40);   /* LSB */   // Least Significant Byte. (Low Byte)
408         outb(LATCH >> 8 , 0x40);       /* MSB */   // Most Significant Byte. (High Byte)
        // Install the handler for timer interrupt. Modify the interrupt mask code to enable system
        // timer generating interrupts. Install the system call gate (interrup 0x80).
409         set_intr_gate(0x20,&timer_interrupt);
410         outb(inb_p(0x21)&~0x01,0x21);
411         set_system_gate(0x80,&system_call);
412 }
413
```

# 5.9  signal.c

## 5.9.1  Function

Program signal.c deals with the all functions processing signals. In the UN*X system, signal is one kind of "soft interrupt" mechanism. A few number of complicated applications will use singals. The signal mechanism provides a way to handle asynchronous events. For example, in the case the user enter combine keys 'ctrl-C' on a terminal keyboard to terminate a running program. This operation will cause the kernel to generate signal SIGINT (signal interrupt), which is to be sent to the foreground process and result in the termination of the process in default situation. Another example is when the process has setup a alarm timer and which was expired, the kernel would send a SIGALRM signal to the process. When an hardware exception occurred, the kernel will also send a corresponding signal to the current process too. In addition, a process may send a signal to the other process. For example, a process send a SIGTERM

signal using function kill() to a process in the same group to terminate it.

The signal processing mechanism is already existed from the old UN*X system, but the way handling sigals in their kernels is not much reliable. The signal may got lost, and the signal is hard to close in the critical code regin. Afterwards, the POSIX provids a reliable way for handling signals using sigaction structure. To keep compatible with these two handle methods, program signal.c provids them at then same time.

Usually, the kernel uses bits in a unsigned long integer (32-bit) to represente different signals for a process. This integer is called signal map in the process descriptor. Thus, there are at most 32 different signals in the system. In the kernel discussed, 22 different signals were defined which includes all the 20 kinds of signals that defined in the POSIX.1 standard. The other two signals are special ones defined by the Linux kernel. They are SIGUNUSED and SIGSTKFLT (stack fault). The former is used to represent all the other signals the system not yet supported. The later is used when the stack encounter error. The name and the difinition of all the 22 different signals are listed in table at the end of the program. You can also refer to the contents of file include/signal.h for the detailed descriptions.

When a signal is received, the process has three different ways to handle it:

1.  Ignore the signal. Most of the signals can be ignored by the process. But the two signals: SIGKILL and SIGSTOP can never be ignored. The reason is that the kernel should provide an affirmatory way for the supervisor to stop or kill any processes. Also, if some of the signals generated by a hardware exception (such as divide-by-zero) are ignored, the behavior of the process is undefined.

2.  Catch the signal. To do this, we should tell the kernel the signal handling function before the it occurred. In the signal function we can do whatever we want to handle the condition, or do nothing to ignore the signal. An example is if the process has created some temporary files, we may want to write a signal-catching function for SIGTERM signal to clean up the temporary files.

3.  Execute the signal's default action. The kernel provides a default action for all of the signals. The generally action of these default signals is to terminate the process.

This program mainly provides the following signal handling functions:

* sys_ssetmask() - System call to set the signal block code of the current process.
* sys_sgetmask() - System call to get the signal block code of the current process.
* sys_signal()   - System call to handle a signal, i.e. the traditional signal handling function signal().
* sys_sigaction() - System call to to change the action taken by a process on receipt of a specific signal, i.e. the reliable signal handling function sigaction().
* do_signal() - The C function invoked in the system calls interrupt.

The action performed by function signal() and sigaction() is similar. They all used to change the signal handler for a indicated signal. But the way the signal() does may cause a signal lost in some particular situations. The additional two functions related to singal handling are send_sig() and tell_father(), and implemented in program exit.c which is the next program to be commented.

At the line 55 in the header file include/signal.h, the prototype of function signal() is declared as follows:

```
    void (*signal(int signr, void (*handler)(int)))(int);
```

The function contains two arguments. The 'signr' is a signal number to catch. The other is a function pointer of the new signal handler which has no return value and contains one integer argument. This integer argument is used by the kernel to send the signal to the signal handler when it occurred.

The signal() installs a new signal handler for the signal with number 'signr'. The signal handler si set to 'handler' which may be a user specified function, or either SIG_IGN or SIG_DFL.

When the specified signal 'signr' arrived, if the corresponding handler is set to SIG_IGN, then the signal is ignored. If the handler is set to SIG_DFL, then the default action associated to the signal performed. Otherwise, if the handler is set to a function 'handler' then first either the handler is reset to SIG_DFL or an implementation-dependent blocking of the signal is performed and next 'handler' is called with argument 'signr'.

On return, the signal() function returns the previous value of the signal handler which is also a function pointer with no return value and contains one integer argument. After the invocation of the new handler, the handler is reset to the default handler SIG_DFL.

At line 45-46 in file include/signal.h, the default handler SIG_DFL and ignoring handler SIG_IGN is defined as a macro as following:

```
    #define SIG_DFL        ((void (*)(int))0)
    #define SIG_IGN        ((void (*)(int))1)
```

They all represent function pointers with no return value same as the second agrument in signal(). The pointer value are 0 and 1 seperately. This two value are function addresses that are impossible to happen in the practical function address. Thus in the signal() function, the kernel can check these two value to determine whether to do the default action or ignore the signal for the process. Of cause the SIGKILL and SIGSTOP can not be ignored. Refer to the L94 - 98 in the signal.c program.

When begin running a program, the kernel will set its signals default handler to SIG_DFL or SIG_IGN. In addition, when forking a child process, this child process will inherit all the signals handling method (signal block mask) from its parent. Thus the setup and handling methods set by parent process are also valid in the child process.

To continuely catch a specified signal, function signal() is usually called in a way like the following example:

```
void sig_handler(int signr)         // Signal handler.
{
    signal(SIGINT, sig_handler);    // Reinstall the signal handler to catch the next signal.
    ...
}

main ()
{
    signal(SIGINT, sig_handler);    // Firstly install the signal handler in the main.
    ...
```

```
}
```

The reason that the signal() is unreliable is that in the case when a signal was happened and entered into its handler, before we re-installed the signal handler again, there may be a new signal occurred in this short period. But the kernel has set the handler to be the default one, thus may cause a signal to be lost.

The function sigaction() adopts sigaction data structure to save the information of the indicated signal. It is a reliable mechanisim to handling signals. We can easily check information and modify the handler of a specified signal. This function is a super-set of the signal() and is declared at L66 in header file include/signal.h like this:

```
int sigaction(int sig, struct sigaction *act, struct sigaction *oldact);
```

This function is used to change the action taken by a process on receipt of a specific signal. The argument 'sig' is the signal number we specified and can be any valid signal except SIGKILL and SIGSTOP. The later two arguments are pointers to the sigaction structure. If 'act' is not null, then the new action for the signal 'sig' is installed from the 'act'. If 'oldact' is non-null, the previous action for the signal is saved in 'oldact'. The sigaction structure is defined as following (in file include/signal.h):

```
48 struct sigaction {
49         void (*sa_handler)(int);       // Signal handler.
50         sigset_t sa_mask;              // Signal mask, used to block specified signal set.
51         int sa_flags;                  // Option flags.
52         void (*sa_restorer)(void);     // Restorer, used internally by the kernel.
53 };
```

'sa_handler' specifies the action to be associated with 'sig' and may be SIG_DFL for the default action, SIG_IGN to ignore this signal, or a pointer to a signal handling function. 'sa_mask' gives a mask of signals which should be blocked during execution of the signal handler. In addition, the signal which triggered the handler will be blocked, unless the SA_NODEFER or SA_NOMASK flags are used. 'sa_flags' specifies a set of flags which modify the behavior of the signal handling process. It is formed by the bitwise OR of zero or more of the flags described in header file include/signal.h, from L35 to L39.

When changing the handling method for a specified signal, if sa_handler is not one of the default handler SIG_IGN or SIG_DFL, then before the 'sa_handler' is to be invoked, the field 'sa_mask' indicates a signal set to be added into the signal block map of the process. If the singal handler returns, the kernel will restore the signal block map for the process. Thus, we can block a set of specified signals before executing the signal handler. In addition, before executing the signal handler, the modified signal block map also contains the current signal 'sig' to prohibit the sending of current signals. Hence, we are guaranteed that whenever we are processing a given signal, another occurrence of that same signal is blocked until we are finished processing the first occurrence. Furthermore, additional occurrences of the same signal are usually not queued. If the signal occurs several times while it is blocked, when we unblock the signal, the signal handling function for that signal will usually be invoked

only one time.

Once we install an action for a given signal, that action remains installed until we explicity change it by calling sigaction again. Unlike the function signals(), POSIX.1 requires that a signal handler remain installed until explicitly changed.

The last field of the sigaction structure and the 'restorer' argument in function sys_signal() is a function pointer, which is provided by the library, for instance Libc, at the compilation period. This function is used to do some cleaning process and restore the return value in register eax by the system call. Refer to the following for detailed descriptions.

Function do_signal() is used to preprocess the signals for a task within the system call interrupt int 0x80. When a process invokes a system call, if the process have received signal(s), then this function will insert the signal handler onto the process's user stack. Thus, when exiting from the system call process, the signal handler will be runned right away, after that, the user program will be continued if the signal handler returns. Refer to Figure 5-8.
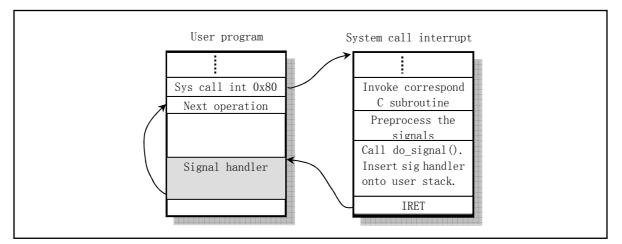


Figure 5-8 The invoking method of the signal handler

Before the insertion of the signal handler onto the user stack, do_signal() will first extends the stack pointer for a few long words ('longs', 7 or 8 long words, see the following program at L106) and, then, filled with arguments in it, as demonstrated in Figure 5-9. Since the code begin from L104 in the following program is hard to understand, here we give them very detailed description.

When the user program (process) invokes system call and begins enter into the kernel code, The kernel stack of the process was filled with some values as shown in Figure 5-9. These values are SS, ESP, EFLAGS and the return address in the user program. When finished handling the corresponding system call function and just before invoking the do_signal() (begin from L118 in program system_call.s), the process's kernel stack contains values as illustrated on the left side in Figure 5-10. This values formed the arguments for the do_signal() function.
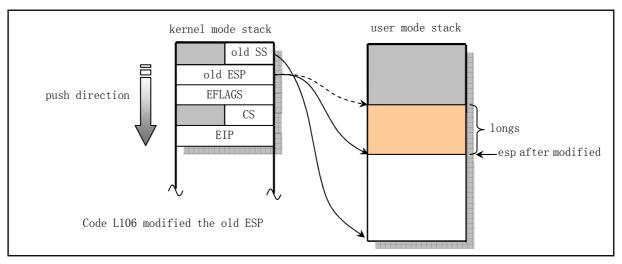
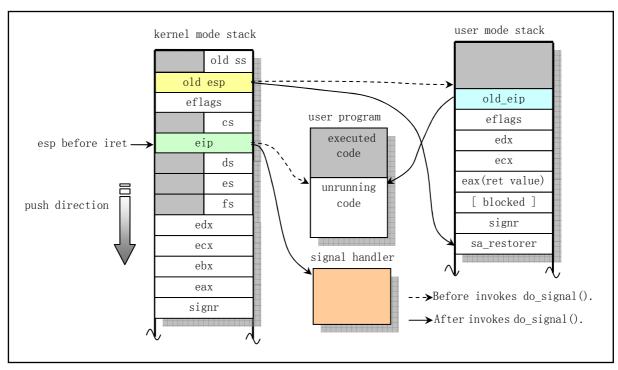Figure 5-9 User stack modification by function do_signal()



Figure 5-10 The detail modification procedure of user stack by do_signal()

After do_signal() deals with the two default signal handlers (SIG_IGN and SIG_DFL), if the user program have defined its own signal handler ('sa_handler'), then do_signal() start prepareing to insert the user signal handler onto the user mode stack at L104. Firstly, it saves the user program return eip in the kernel mode stack to be as old_eip and replace the eip with the pointer of the user defined signal handler 'sa_handler'. Then, it expands the current using space of the user mode stack by decreasing 'old esp' in 'longs' size. Lastly, it copies a few register values on the kernel mode stack onto the user mode stack as shown in the right side of Figure 5-10.

The total number of parameters putted on the user stack are 7 or 8 long words. We now describe their meanings below.

'old_eip' is the orignal return instruction pointer of the user program. It was saved before

the eip be replaced with the user signal handler on the kernel mode stack. 'eflags', 'edx', 'ecx' are the values when user program invokes the system call and are basically the arguments for the system call. When end of the system call, we need to restore these values in the relating registers. 'eax' keeps the system call return code. If the signal allow the reception of itself, then we also need to store the process signal block mask 'blocked' onto the user stack. This is reason why the 'longs' size is one long words variable. The next value is signal number 'signr'. The last one is a signal action restore function pointer 'sa_restorer'. This function is not provided by the user. The user only defines the signal number 'signr' and the signal handler 'handler' for the function signal().

The following is a simple example to setup a signal handler for the signal SIGINT. In the default case, this signal is generated when we press Ctrl-C keys.

```c
#include <signal.h>
#include <stdio.h>
#include <unistd.h>
// Signal handler.
void handler(int sig)                     // Signal processing function.
{
    printf("The signal is %d\n", sig);
    (void) signal(SIGINT, SIG_DFL);       // Restore the default signal handler.
}

int main()
{
    (void) signal(SIGINT, handler);       // Install the user defined signal handler.
    while (1) {
        printf("Signal test.\n");
        sleep(1);                         // Wait one second.
    }
}
```

The signal processing function handler() will be invoked when the process receives a SIGINT signal. The handler() first shows a message and then restores the default signal handler for signal SIGINT. Thus when we press the combining keys CTRL-C the second times, this program will be terminated.

Now, lets discuss the function sa_restorer. We cannot find its difinition in the kernel sources, so, where does it come from? As a matter of fact, it was definied in the library file libc.a as following:

```asm
.globl ____sig_restore
.globl ____masksig_restore
# If there is no 'blocked' then use procedure restorer
____sig_restore:
    addl $4,%esp        # get rid of sinal value 'signr'.
    popl %eax           # restore the return value of the system call.
    popl %ecx           # restore the original register values of the user program.
    popl %edx
```

```
    popfl                   # restore the EFLAGS of user program.
    ret
# If there is a 'blocked', then we use the following restorer procedure.
# 'blocked' is used by ssetmask.
____masksig_restore:
    addl $4,%esp            # ignore the signal value 'signr'.
    call ____ssetmask       # Setup signal mask 'old blocking'.
    addl $4,%esp            # discard 'blocked'
    popl %eax
    popl %ecx
    popl %edx
    popfl
    ret
```

The main purpose of the above function is to restore the returnning value, a fews register values and discard the signal value on the stack when finished executing the system calll by the user program at the end of the processing the signal handler. When compiling and linking the user defined signal handling function, the compiler will insert the above function into the user's program and will do the above job for the program as if the kernel return directly from the system call without run the user defined signal handler.

Now we give a sum description of the signal.c program. After finished executing function do_signal(), program system_call.s pops out all values bellow the value of eip on the kernel state stack. After running the instruction 'iret', the CPU pops out cs:eip, eflags and ss:esp and restored to the user state and continue running the user program. Since the eip has been replaced with the signal handler pointer, the CPU executes the user defined signal handler right away. After the executing of signal handler, the CPU passes the control to the function sa_restorer. This function does some clearing and restoring jobs for the user program and return and continue to run the user program.

## 5.9.2 Comments

Program 5-7 linux/kernel/signal.c

```
 1 /*
 2  *  linux/kernel/signal.c
 3  *
 4  *  (C) 1991  Linus Torvalds
 5  */
 6
 7 #include <linux/sched.h>
 8 #include <linux/kernel.h>
 9 #include <asm/segment.h>
10
11 #include <signal.h>
12
13 volatile void do_exit(int error_code);
14
    // Get the mask of blocked signals.
15 int sys_sgetmask()
```

```
16 {
17         return current->blocked;
18 }
19
   // Set the new signal blocking mask. SIGKILL can not be blocked. ret the old blocking mask.
20 int sys_ssetmask(int newmask)
21 {
22         int old=current->blocked;
23
24         current->blocked = newmask & ~(1<<(SIGKILL-1));
25         return old;
26 }
27
   // Copy the data of sigaction to the address 'to' in the user space (fs segment).
28 static inline void save_old(char * from,char * to)
29 {
30         int i;
31
   // Verify that there is enough memory space at the address 'to' to hold the data.
32         verify_area(to, sizeof(struct sigaction));
33         for (i=0 ; i< sizeof(struct sigaction) ; i++) {
34                 put_fs_byte(*from,to);          // Copy to segment fs of user space.
35                 from++;                         // put_fs_byte() is in include/asm/segment.h
36                 to++;
37         }
38 }
39
   // Get the sigaction data from user space (fs segment).
40 static inline void get_new(char * from,char * to)
41 {
42         int i;
43
44         for (i=0 ; i< sizeof(struct sigaction) ; i++)
45                 *(to++) = get_fs_byte(from++);
46 }
47
   // The signal() system call. It installs a new signal handler for the signal with number
   // 'signum'. The signal handler is set to 'handler' which may be a user specified function
   // or either SIG_IGN or SIG_DFL.
   // Arguments: signum -- the specified signal;
   //            handler -- the signal handler;
   //            restorer -- the restoring function pointer.
   // The restoring function is provided by library file libc.a. It is used to do the ending
   // process for the user program after executing the signal handler.
   // This system call returns the orignal signal handler.
48 int sys_signal(int signum, long handler, long restorer)
49 {
50         struct sigaction tmp;
51
   // the signal number must in the range of 1 to 32, and cannot be signal SIGKILL.
52         if (signum<1 || signum>32 || signum==SIGKILL)
53                 return -1;
   // This signal uses the handler only once and then restore to the default handler.
```

```
     // The signal can received in its handler.
54          tmp.sa_handler = (void (*)(int)) handler;      // The specified signal handler.
55          tmp.sa_mask = 0;                               // signal mask when running.
56          tmp.sa_flags = SA_ONESHOT | SA_NOMASK;
57          tmp.sa_restorer = (void (*)(void)) restorer;
58          handler = (long) current->sigaction[signum-1].sa_handler;
59          current->sigaction[signum-1] = tmp;
60          return handler;
61  }
62
     // The sigaction() system call is used to change the action taken by a process on receipt
     // of a specific signal. 'signum' specifies the signal and can be any valid signal except
     // SIGKILL (and SIGSTOP). If 'action' is non-null, the new action for signal 'signum' is
     // installed from 'action'. If 'oldaction' is non-null, the previous action is saved in it.
     // The function returns 0 on success and -1 on error.
63  int sys_sigaction(int signum, const struct sigaction * action,
64          struct sigaction * oldaction)
65  {
66          struct sigaction tmp;
67
     // The valid signal numer is in the range of 1 to 32 and cannot be signal SIGKILL.
68          if (signum<1 || signum>32 || signum==SIGKILL)
69                  return -1;
     // Get the user defined sigaction data from 'action' in user space to the current process.
70          tmp = current->sigaction[signum-1];
71          get_new((char *) action,
72                  (char *) (signum-1+current->sigaction));
     // If 'oldaction' is not null, then save the orignal sigaction data for the user.
73          if (oldaction)
74                  save_old((char *) &tmp,(char *) oldaction);
     // If we permit to receive the signal in the signal's handler, then clear the signal mask,
     // else set to mask the specified signal 'signum'.
75          if (current->sigaction[signum-1].sa_flags & SA_NOMASK)
76                  current->sigaction[signum-1].sa_mask = 0;
77          else
78                  current->sigaction[signum-1].sa_mask |= (1<<(signum-1));
79          return 0;
80  }
81
     // The signal preprocessing function for the system call interrupt procedure invoked at
     // line 119 in kernel/system_call.s file. This function is used to insert the signal
     // handler into the user state stack for executing the signal handler right after return
     // from the system call.
82  void do_signal(long signr,long eax, long ebx, long ecx, long edx,
83          long fs, long es, long ds,
84          long eip, long cs, long eflags,
85          unsigned long * esp, long ss)
86  {
87          unsigned long sa_handler;
88          long old_eip=eip;
89          struct sigaction * sa = current->sigaction + signr - 1; //current->sigaction[signu-1]。
90          int longs;
91          unsigned long * tmp_esp;
```

```
92
93          sa_handler = (unsigned long) sa->sa_handler;
    // If the signal handler is SIG_IGN (ignore the signal), then return; if the signal handler
    // is SIG_DFL (the default handler), then, if the signal is SIGCHILD then return, else we
    // terminate the program.
    // The handler SIG_IGN is defined as 1, SIG_DFL is defined as 0. Refer to line 45, 46 in
    // file include/signal.h.
94          if (sa_handler==1)
95                  return;
96          if (!sa_handler) {
97                  if (signr==SIGCHLD)
98                          return;
99                  else
100                         do_exit(1<<(signr-1));   // This should be do_exit(1<<(signr)).
101         }
    // If the handler need to run once, then clear the handler for this signal.
102         if (sa->sa_flags & SA_ONESHOT)
103                 sa->sa_handler = NULL;
    // The following code is used to insert the signal handler onto user stack, and at the same
    // time pushes the restore function pointer 'sa_restorer', signal number 'signr', process
    // mask code (if SA_NOMASK isn't set), register eax, ecx, edx, eflags and the user return
    // address on to the user stack. Thus, the CPU will run the user specified signal handler
    // first and then continue to execute the user program after the system call.

    // Change the value of eip on the kernel state stack to the pointer of the signal handler.
104         *(&eip) = sa_handler;
    // If it allows the signal handler to receive the own signal, then we need to push the
    // signal blocking mask too. Notice that the result of 'longs' is one of (7*4) and (8*4).
105         longs = (sa->sa_flags & SA_NOMASK)?7:8;
    // To store the arguments of signal handler, we need to extend the space for them on the
    // user state stack and varify the space (to allocate new page if not enough).
106         *(&esp) -= longs;
107         verify_area(esp,longs*4);
    // In the user state stack, from the bottom to top, we store sa_restorer, signal number signr,
    // signal blocking mask, register valuse eax, ecx, edx, eflags and the user program pointer.
108         tmp_esp=esp;
109         put_fs_long((long) sa->sa_restorer,tmp_esp++);
110         put_fs_long(signr,tmp_esp++);
111         if (!(sa->sa_flags & SA_NOMASK))
112                 put_fs_long(current->blocked,tmp_esp++);
113         put_fs_long(eax,tmp_esp++);
114         put_fs_long(ecx,tmp_esp++);
115         put_fs_long(edx,tmp_esp++);
116         put_fs_long(eflags,tmp_esp++);
117         put_fs_long(old_eip,tmp_esp++);
118         current->blocked |= sa->sa_mask;            // Add signal mask to blocked in PCB.
119 }
120
```

## 5.9.3  Information

### 5.9.3.1  Description of process signals

The signal in the process is a simple message used to communicate between processes. It

is a simple integer value as shown in Table 5-3. For example, while a child process stop running or terminates, it will generate the signal SIGCHILD and send to the father to notify the current status of the child process.

To handling the signal received, a process has two kinds of method. One way is to do nothing and the kernel takes the default action for the signal as listed in the right side of the table. The other method is to write a function or handler for the signal and to handle the signal by the program itself.

Table 5-3 Process signals

| Value | Name | Description | Default Action |
|-------|------|-------------|----------------|
| 1 | SIGHUP | (Hangup) The kernel generates this signal when the process has no controlling terminal or when the modem is hangup.Since daemons doesn't have controlling terminal, they often use this signal to require rereading configuration file. | (Abort) Hang up the controlling terminal or process. |
| 2 | SIGINT | (Interrupt) Interrupt from keyboard. Usually the terminal driver combines this signal with key stroke 'ˆC'. | (Abort) Terminate the process. |
| 3 | SIGQUIT | (Quit) Interrupt from keyboard. Usually the terminal driver combines this signal with key stroke ˆ\. | (Dump) Terminate the process and dump core. |
| 4 | SIGILL | (Illegal Instruction) Program error or executes a illegal instruction. | (Dump) Terminate the process and dump core. |
| 5 | SIGTRAP | (Breakpoint/Trace Trap) Breakpoint for debugging. | (Dump) Terminate the process and dump core. |
| 6 | SIGABRT | (Abort) Abort the execution, abnormal termination. | (Dump) Terminate the process and dump core. |
| 6 | SIGIOT | (IO Trap) Hardware error, equivalent to SIGABRT. | (Dump) Terminate the process and dump core. |
| 7 | SIGUNUSED | (Unused) Not used, bad system call. | (Abort) Terminate the process. |
| 8 | SIGFPE | (Floating Point Exception) FPU error. | (Dump) Terminate the process and dump core. |
| 9 | SIGKILL | (Kill) Terminate the program. This signal can not be captured or ignored. If you want to terminate a program immediately, send this signal. | (Abort) Terminate the process. |
| 10 | SIGUSR1 | (User defined Signal 1) Available to processes. | (Abort) Terminate the process. |
| 11 | SIGSEGV | (Segmentation Violation) Invalid memory reference. | (Dump) Terminate the process and dump core. |
| 12 | SIGUSR2 | (User defined Signal 2) Available to processes for IPC or other purposes. | (Abort) Terminate the process. |
| 13 | SIGPIPE | (Pipe) Write to pipe with no readers. | (Abort) Terminate the process. |
| 14 | SIGALRM | (Alarm) Generated when the alarm timer expired. | (Abort) Terminate the process. |

| 15 | SIGTERM | (Terminate) Require to terminate a process friendly. It is the default signal for kill(). This signal can be captured. | (Abort) Terminate the process. |
|----|---------|------|------|
| 16 | SIGSTKFLT | (Stack fault on coprocessor) | (Abort) Terminate the process. |
| 17 | SIGCHLD | (Child) Send to father process and tell him the child process has stopped or terminated. | (Ignore) Ignore this signal. |
| 18 | SIGCONT | (Continue) Resume the stopped process to run state. | (Continue) Resume execution if stopped. |
| 19 | SIGSTOP | (Stop) Stop process execution. This signal can not be captured or ignored. | (Stop) Stop the process. |
| 20 | SIGTSTP | (Terminal Stop) Stop the process issued from terminal. This signal can be captured or ignored. | (Stop) Stop the process. |
| 21 | SIGTTIN | (TTY Input on Background) The background process try to do input operation from a terminal. This process will be stopped until receiving a SIGCONT signal. | (Stop) Stop the process. |
| 22 | SIGTTOU | (TTY Output on Background) The background process try to do output operation to a terminal. This process will be stopped until receiving a SIGCONT signal. | (Stop) Stop the process. |

# 5.10 exit.c

## 5.10.1 Function

This program primarily describes the operation of the termination or exit of a process. The actions involve process' resources release, session (or process group) termination and system calls like killing a process, terminating a process and hang up a process. It also includes signal sending functions such as send_sig() and tell_father().

The release() function deletes the process pointer in the task array, release the memory page relating to the process task structure and finally rescheduling the task.

The kill_session() function kills all the process in the session by sending SIGHUP signal.

The kill() system call is used to send any valid signal to a process. Based on the argument 'pid', this system call can send signal to deferent process or process group. The comments in the program listed bellow has presented all the cases.

The program exit handling function do_exit() is invoked in the interrupt procedure of the exit system call. Firstly, it frees the memory pages occupied by the code and data segments of the current process. If the current process owns the child processes, the father of them will be set to init process ( process 1). If one of the child process has already be in zombie state, a SIGCHLD signal will be send to init process for this child process. Secondly, it closes all the files, releases the terminal device and coprocessor device. If the current process is the header of the process group, then all process in the group will be terminated. Thirdly, the current process is set to zombie state, the exit code is set and a signal SIGCHLD is sent

to its father process. Lastly, it does task rescheduling.

The system call waitpid() is used to hang up the current process until the child process specified by the pid value is terminated or a terminating signal is received, or a signal handler is need to be invoked. If the child process specified by pid has already exit, waitpid() returns immediately. All resource used by the child process will be released.

## 5.10.2 Comments

Program 5-8 linux/kernel/exit.c

```
1  /*
2   *  linux/kernel/exit.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  #include <errno.h>
8  #include <signal.h>
9  #include <sys/wait.h>
10
11 #include <linux/sched.h>
12 #include <linux/kernel.h>
13 #include <linux/tty.h>
14 #include <asm/segment.h>
15
16 int sys_pause(void);      // System call causing the process to sleep. (L144 in sched.c)
17 int sys_close(int fd);    // System call closing the file. (L192 in fs/open.c)
18
   //// Free the task slot and release the memory page occupied by the task p.
   // Argument p is a task structure pointer. This function scan the task table task[] to
   // find the specified task p. If found, it first free the task slot of p, then release the
   // memory page occupied by the task structure and its kernel stack space, and lastly doing
   // reschedule. If the task is not found, then the kernel goes panic.
   // This function is invoked in the function sys_kill() and sys_waitpid().
19 void release(struct task_struct * p)
20 {
21        int i;
22
23        if (!p)
24                return;
25        for (i=1 ; i<NR_TASKS ; i++)        // Scan task table for the task p.
26                if (task[i]==p) {
27                        task[i]=NULL;       // Free the task slot.
28                        free_page((long)p); // Free the memory page.
29                        schedule();         // Rescheduling (Really needed?).
30                        return;
31                }
32        panic("trying to release non-existent task");
33 }
34
   //// Send signal 'sig' to the specified task 'p'.
   // Arguments 'sig' is the signal number; 'p' is the specified task structure pointer;
```

```
     // 'priv' is a signal sending privilege, when set, a signal can be sent without care
     // about the user property or level of the sending process. This function checks the
     // correctness of the arguments and send the signal to the specified task if the constraint
     // condition ('priv') is satisfied.
35   static inline int send_sig(long sig, struct task_struct * p, int priv)
36   {
37           if (!p || sig<1 || sig>32)
38                   return -EINVAL;
     // If the signal sending privilege is set, or the euid (effective uid) of the current
     // process is the euid of the specified process, or the current process is the supervisor,
     // then the signal 'sig' is sent to the process 'p'. That is the corresponding bit in the
     // task's field 'signal' is set. The suser() is defined as (current->euid==0), and used to check
     // if the process is owned by the root user.
39           if (priv || (current->euid==p->euid) || suser())
40                   p->signal |= (1<<(sig-1));
41           else
42                   return -EPERM;
43           return 0;
44   }
45
     //// Terminate the current session.
46   static void kill_session(void)
47   {
     // The task pointer *p is first set to point to the last entry of the task table.
48           struct task_struct **p = NR_TASKS + task;
49
     // Scan the task pointer table. For the all tasks, except task 0, if the task's session
     // number equals to that of the current task, then sends signal SIGHUP to the task.
50           while (--p > &FIRST_TASK) {
51                   if (*p && (*p)->session == current->session)
52                           (*p)->signal |= 1<<(SIGHUP-1);
53           }
54   }
55
56   /*
57    * XXX need to check permissions needed to send signals to process
58    * groups, etc. etc.  kill() permissions semantics are tricky!
59    */
     //// The system call kill() can send any kind of signal to any process or process group,
     // not only can just kill a process.
     // Argument 'pid' is the specified processs's pid; 'sig' is the specified signal number.
     // If pid is positive, then signal sig is sent to the process with pid.
     // If pid equals 0, then sig is sent to every process in the process group of the current
     // process.
     // If pid equals -1, then sig is sent to every process except for process 1 (init).
     // If pid is less than -1, then sig is sent to every process in the process group -pid.
     // If sig is 0, then no signal is sent, but error checking is still performed.
     // The system call function scans the task table and sends the signal to the tasks that
     // satisfied the conditions. If pid equals 0, it means that the current process is the
     // leader of a process group. Thus, the signal needs to be sent to every process in the
     // process group of the current process.
60   int sys_kill(int pid, int sig)
61   {
```

```
62          struct task_struct **p = NR_TASKS + task;
63          int err, retval = 0;
64
65          if (!pid) while (--p > &FIRST_TASK) {
66                  if (*p && (*p)->pgrp == current->pid)
67                          if (err=send_sig(sig,*p,1))          // Force to send the signal.
68                                  retval = err;
69          } else if (pid>0) while (--p > &FIRST_TASK) {
70                  if (*p && (*p)->pid == pid)
71                          if (err=send_sig(sig,*p,0))
72                                  retval = err;
73          } else if (pid == -1) while (--p > &FIRST_TASK)
74                  if (err = send_sig(sig,*p,0))
75                          retval = err;
76          else while (--p > &FIRST_TASK)
77                  if (*p && (*p)->pgrp == -pid)
78                          if (err = send_sig(sig,*p,0))
79                                  retval = err;
80          return retval;
81  }
82
```
//// Notifying the father process. Send signal SIGCHLD to the process pid. The child
// process will be terminated or stopped in the default case.
// If the father process is not found, then the operated process need to be freeed by
// itself. But for the requirement of POSIX.1, if the father process terminated previously,
// its child processes should be adopted by the init process (process 1).
```
83  static void tell_father(int pid)
84  {
85          int i;
86
87          if (pid)
```
// Scan the process table to find the specified process of pid. Send the signal SIGCHLD to
// the process pid, the child process will be stopped or terminated.
```
88                  for (i=0;i<NR_TASKS;i++) {
89                          if (!task[i])
90                                  continue;
91                          if (task[i]->pid != pid)
92                                  continue;
93                          task[i]->signal |= (1<<(SIGCHLD-1));
94                          return;
95                  }
96  /* if we don't find any fathers, we just release ourselves */
97  /* This is not really OK. Must change it to make father 1 */
98          printk("BAD BAD - no father found\n\r");
99          release(current);                       // Release the process itself if no father found.
100 }
101
```
//// Program exits handling function, called in the system call sys_exit() at L137.
// Argument 'code' is a error code.
```
102 int do_exit(long code)
103 {
104         int i;
105
```

```
      // Free the page table entries corresponding to the memory pages occupied by the code and
      // data segment of the current process. 'free_page_tables()' is at L105 in mm/memory.c
106         free_page_tables(get_base(current->ldt[1]),get_limit(0x0f));
107         free_page_tables(get_base(current->ldt[2]),get_limit(0x17));
      // If the current process has got child processes, then changes the father process of the
      // child processes to be process init (process 1). If a child is in the state of zombie,
      // then sends a signal SIGCHLD to the init process.
108         for (i=0 ; i<NR_TASKS ; i++)
109               if (task[i] && task[i]->father == current->pid) {
110                     task[i]->father = 1;
111                     if (task[i]->state == TASK_ZOMBIE)
112           /* assumption task[1] is always init */
113                           (void) send_sig(SIGCHLD, task[1], 1);
114               }
      // Close all files opened.
115         for (i=0 ; i<NR_OPEN ; i++)
116               if (current->filp[i])
117                     sys_close(i);
      // Synchronizing the working directory, root directory and the I node of the executing
      // program, then release them.
118         iput(current->pwd);
119         current->pwd=NULL;
120         iput(current->root);
121         current->root=NULL;
122         iput(current->executable);
123         current->executable=NULL;
      // If the current process is the session leader and has a controlling terminal, then
      // release it.
124         if (current->leader && current->tty >= 0)
125               tty_table[current->tty].pgrp = 0;
      // If the current process once used coprocessor, then clear the used flag.
126         if (last_task_used_math == current)
127               last_task_used_math = NULL;
      // If the current process is the session leader, then kills all the processes in the session.
128         if (current->leader)
129               kill_session();
      // Change the state of the current process to be ZOMBIE, this means that the current process
      // has released the resources, and saved the exit code to be read for the father process.
130         current->state = TASK_ZOMBIE;
131         current->exit_code = code;
      // Send signal SIGCHLD to the father process.
132         tell_father(current->father);
      // Rescheduling, the father process can deal with the other matters concerned.
133         schedule();
134         return (-1);    /* just to suppress warnings */
135 }
136
      //// System call exit(). Terminating the process.
137 int sys_exit(int error_code)
138 {
139         return do_exit((error_code&0xff)<<8);
140 }
141
```

```
      // System call waitpid(). It suspends the execution of the current process until a child
      // process as specified by the pid argument has exited, or until a signal is delivered whos
      // action is to terminated the current process or to call a signal handling function. If a
      // child process as requested by pid has already exited by the time of the call, the function
      // returns immediately. Any system resources used by the child process aare freed.
      // If pid great than 0, then which means to wait for the child process whose process Id
      // equals to the vaule of pid.
      // If pid equals 0, which means to wait for any child process whose process group ID is equal
      // to that of the calling process.
      // If pid is less than -1, which means to wait for any child process whose process group ID
      // is equal to the absolute value of pid.
      // If pid is equal -1, which means to wait for any child process.
      // The value of 'options' is an OR of zero or more of the following constants:
      // WNOHANG, which means to return immediately if no child has exited.
      // WUNTRACED, which means to also return for children which are stopped, and whose status
      // has not been reported.
      // If 'stat_addr' is not NULL, waitpid store status information in the location pointed
      // by 'stat_addr'.
142 int sys_waitpid(pid_t pid,unsigned long * stat_addr, int options)
143 {
144         int flag, code;
145         struct task_struct ** p;
146
147         verify_area(stat_addr,4);
148 repeat:
149         flag=0;
      // Scan the task table begins from the last entry, skip the empty slot, current process
      // and the child processes not belong to the current process.
150         for(p = &LAST_TASK ; p > &FIRST_TASK ; --p) {
151             if (!*p || *p == current)            // skip the empty and current slot.
152                 continue;
153             if ((*p)->father != current->pid)    // skip the non current's children.
154                 continue;
      // The selected process p from the scanning is definitively the child of current process.
      // If the waited pid great than 0, but it is not equal to the p's pid, which means the p
      // is one of the other child of the current process, then we skip it.
155             if (pid>0) {
156                 if ((*p)->pid != pid)
157                     continue;
      // Otherwise, if the specified waited pid is equal 0, which means we are waiting for any
      // child whose process group ID is equal to the current process. In this case, if the
      // scanning process's group ID is not equal to that of the current process, then we skip.
158             } else if (!pid) {
159                 if ((*p)->pgrp != current->pgrp)
160                     continue;
      // Otherwise, if the pid is less than -1, which means we are waiting for any child process
      // whose process group ID is equal to the absolute value of pid. In this case, if the
      // scanning process's group ID is not equal to the absolute value of the pid, then we skip.
161             } else if (pid != -1) {
162                 if ((*p)->pgrp != -pid)
163                     continue;
164             }
      // If the above three judgements are not correct, it means the current process is waiting
```

```
         // for any of its child process. This is the case with pid equals -1.
         // At this moment, the selected process p is exactly the waited process. Now, we handle
         // this child process based upon its status.
165                  switch ((*p)->state) {
         // In the case the child is in STOPPED status, if the flag WUNTRACED is not set, which means
         // we needn't return immediately, then we continue to scan the other processes. If the flag
         // is set, then we save the status information 0x7f at the location of '*stat_addr' and
         // return with the value of child's pid. Here the status value of 0x7f will make the macro
         // WIFSTOPPED() to be true. Refer to L14 in include/sys/wait.h
166                      case TASK_STOPPED:
167                          if (!(options & WUNTRACED))
168                              continue;
169                          put_fs_long(0x7f, stat_addr);
170                          return (*p)->pid;
         // In the case the child is in ZOMBIE status, we first add up its running time in user
         // mode and kernel mode to the current process (the father process) respectively, then we
         // obtain the child's pid and exit code and release the child. Finally we return the exit
         // code and child's pid.
171                      case TASK_ZOMBIE:
172                          current->cutime += (*p)->utime;
173                          current->cstime += (*p)->stime;
174                          flag = (*p)->pid;                   // Save child's pid.
175                          code = (*p)->exit_code;             // Get child's exit code.
176                          release(*p);                        // Release the child.
177                          put_fs_long(code, stat_addr);       // Store the exit code.
178                          return flag;                        // Return the child's pid.
         // If the child's status is not STOPPED and ZOMBIE, then we temporary set the flag = 1, which
         // means we have got one child matched condition but in the state of running or sleeping, and
         // will be used bellow.
179                      default:
180                          flag=1;
181                          continue;
182                  }
183          }
         // When finished scanning, if the flag is set, which means there is at least one child
         // that matches the condition (not stopped or zombie). At this moment, if the option
         // WNOHANG is set, which means to return immediately if no child has exited, then we
         // return right away, otherwise we set the status of current process to be INTERRUPTIBLE
         // and do reschedule.
184          if (flag) {
185              if (options & WNOHANG)          // Return immediately if options = WNOHANG
186                  return 0;
187              current->state=TASK_INTERRUPTIBLE;
188              schedule();
         // When re-executing this process, if there is no signal received (except SIGCHILD), then
         // repeat the operation, otherwise return with error code.
189              if (!(current->signal &= ~(1<<(SIGCHLD-1))))
190                  goto repeat;
191              else
192                  return -EINTR;
193          }
         // If we couldn't find a matched child, then return with error code.
194          return -ECHILD;
```

195 }
196

# 5.11 fork.c

## 5.11.1 Function

fork system call is used to create new (child) process. In the Linux system, all processes are child process of init process (process 1) except process 0, which is created manually in the system initialization period. fork system call is implemented in the file system_call.s as a label name of sys_fork. The fork.c program contains assistant functions for the sys_fork.

The two main functions are find_empty_process() and copy_process(). In file system_call.s, fork system call first invokes find_empty_process() to find and get a empty slot in the process table, then invokes copy_process() to duplicate the task structure of the current (father) process. The other two important functions are memory area verification function verify_area() and memory copy function copy_mem().

The copy_process() function is used to create and duplicate the code, data and environment information of the current process for new child process. In the procedure of copying, the main job is to set up the parameters in the new process's task structure. Firstly, the function allocates one page in the main memory area for storing the new task structure, then it duplicates the information in the current process's task structure to be used as a template. Secondly, it modifies the contents of the new task structure, sets the current process to be the father of the new process, clears signal map and resets the accounting fields. Thirdly, it sets up the contents of the TSS based on that of the current process.

Since the new child needs to return 0, we should set tss.eax = 0. The new child's kernel stack pointer tss.esp0 should be set to the location at the top of the new allocated memory page, which is also occupied by the task structure begin from the bottom of the page. The new child's kernel stack segment register tss.ss0 is set to equal to the segment selector of the kernel data segment. The tss.ldt is set to the index value of the relating ldt entry in the GDT. If the current process has been using the coprocessor, then the whole context of the coprocessor should be saved to the tss.i387 structure in the new task structure.

After the above actions, the copy_mem() function is called to set up the base and size limits of the new task, and duplicates the pages tables of the current process. If there are files opened in the current process (the father process), these files are also be kept opening in the new process. So the coresponding file open counts need to be increased by 1 respectively. Hereafter, the TSS and LDT descriptor entries for the new process is to be installed in the GDT. The base addresses of their descriptors are pointed to the tss and ldt fields in the new task structure.

Finally, the status of the new process is set to be runnable, and the pid of the new process is returned to the current process.

Figure 5-11 demonstrates the adjusting range and start address of the memory for the

verify_area() function. Since the write_verify() function needs to operate with unit in pages
(4096 bytes per page), the orignal start address of the verifying memory space should be adjusted
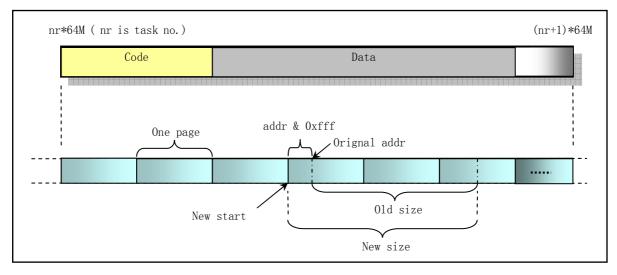to be right at the pages boundary, and the verifying memory size is also adjusted accordingly.



Figure 5-11 Adjusting the range and start address of the verifying memory.

## 5.11.2 Comments

Program 5-9 linux/kernel/fork.c

```
1  /*
2   *  linux/kernel/fork.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  /*
8   *  'fork.c' contains the help-routines for the 'fork' system call
9   * (see also system_call.s), and some misc functions ('verify_area').
10  * Fork is rather simple, once you get the hang of it, but the memory
11  * management can be a bitch. See 'mm/mm.c': 'copy_page_tables()'
12  */
13 #include <errno.h>
14
15 #include <linux/sched.h>
16 #include <linux/kernel.h>
17 #include <asm/segment.h>
18 #include <asm/system.h>
19
20 extern void write_verify(unsigned long address);   // At L261 in mm/memory.c
21
22 long last_pid=0;               // The newest pid number, generated by get_empty_process().
23
   // Function of verifying memory space before write operation.
   // Checking the memory space ranging from 'addr' to 'addr + size' before write operation
   // using page unit for the current process. Since the checking action is based on pages,
```

```
   // it needs to find the start address 'start'of the page that contains the 'addr', and
   // convert it to the CPU 4GB linear space by adding the base address of the process data
   // segment. Afterwards, the write_verify() function is called in circle to do the real
   // checking operation. If the page that will be written is readonly, then share verification
   // and page copy operation are performend (Copy on Write).
24 void verify_area(void * addr,int size)
25 {
26         unsigned long start;
27
28         start = (unsigned long) addr;
   // Adjust the beginning address 'addr' to the left boundary of the page that contains
   // the 'addr', and resizing the verifying size accordingly.
   // 'start & 0xfff' is used to obtains the offset in a page. By by adding the orignal size,
   // the result 'size' is a adjusted size for the start address begins from a page boundary.
   // The new start address is adjusted as 'start & 0xfffff000'.
29         size += start & 0xfff;
30         start &= 0xfffff000;
   // Since the write_verify() function is operated in CPU 4GB linear space, the address 'start'
   // needs first to be converted into it by adding the base address of the process.
31         start += get_base(current->ldt[2]);
32         while (size>0) {
33                 size -= 4096;
   // Memory write verifying. If a page can not be written, then duplicates the page.
34                 write_verify(start);              // At L261 in mm/memory.c
35                 start += 4096;
36         }
37 }
38
   // Setup the base addresses and size limits of the code and data segments, copy page tables.
   // Arguments nr is the new task number; p is the pointer of the new task structure.
39 int copy_mem(int nr,struct task_struct * p)
40 {
41         unsigned long old_data_base,new_data_base,data_limit;
42         unsigned long old_code_base,new_code_base,code_limit;
43
44         code_limit=get_limit(0x0f);    // Get size limit in the code seg descriptor in LDT.
45         data_limit=get_limit(0x17);    // Get size limit in the data seg descriptor in LDT.
46         old_code_base = get_base(current->ldt[1]);    // Get code base of current process.
47         old_data_base = get_base(current->ldt[2]);    // Get data base of current process.
48         if (old_data_base != old_code_base)
49                 panic("We don't support separate I&D");   // Instruction & Data area.
50         if (data_limit < code_limit)
51                 panic("Bad data_limit");
52         new_data_base = new_code_base = nr * 0x4000000;   // New base = Task nr * 64MB.
53         p->start_code = new_code_base;
54         set_base(p->ldt[1],new_code_base);                // Setup base address for new task.
55         set_base(p->ldt[2],new_data_base);
   // Copy page table entries from the current process space to new task space. If an error
   // occurred during copying, then release them all and return error code.
56         if (copy_page_tables(old_data_base,new_data_base,data_limit)) {
57                 free_page_tables(new_data_base,data_limit);
58                 return -ENOMEM;
59         }
```

```
60              return 0;
61 }
62
63 /*
64  *  Ok, this is the main fork-routine. It copies the system process
65  * information (task[nr]) and sets up the necessary registers. It
66  * also copies the data segment in it's entirety.
67  */
   // Copy and manipulate the new task structure from the current process.
   // Argument 'nr' is the new task number, which is the return value of find_empty_process().
   // 'none' is the return address pushed onto the stack when invoking sys_call_table in file
   // system_call.s
68 int copy_process(int nr,long ebp,long edi,long esi,long gs,long none,
69                  long ebx,long ecx,long edx,
70                  long fs,long es,long ds,
71                  long eip,long cs,long eflags,long esp,long ss)
72 {
73         struct task_struct *p;
74         int i;
75         struct file *f;
76
   // Get one free page from main memory area for the new task structure and its kernel stack
   // space. Put the new task structure pointer into task table in slot nr. Copy the contents
   // of current process's task structure into the new task structure.
77         p = (struct task_struct *) get_free_page();
78         if (!p)
79                 return -EAGAIN;
80         task[nr] = p;
81         *p = *current;                  /* NOTE! this doesn't copy the supervisor stack */
   // We first set the status of the new task to be uninterruptible wait state. setup the
   // new process pid to be 'last_pid', which is a global variable evaluated in the function
   // find_empty_process().
82         p->state = TASK_UNINTERRUPTIBLE;
83         p->pid = last_pid;              // New process ID.
84         p->father = current->pid;       // Set the fahter.
85         p->counter = p->priority;
86         p->signal = 0;                  // Reset the signal map. No signals now.
87         p->alarm = 0;                   // Alarm timer (ticks).
88         p->leader = 0;                  /* process leadership doesn't inherit */
89         p->utime = p->stime = 0;        // The accounting values.
90         p->cutime = p->cstime = 0;      // Child process accounting values.
91         p->start_time = jiffies;        // Current system time.
   // The following code initializes the TSS for the new task. Since the kernel allocate
   // one memory page for the new task, thus, the tss.esp0 right points to the top of the
   // page. ss0:esp0 is used as the kernel mode stack for the new task.
92         p->tss.back_link = 0;
93         p->tss.esp0 = PAGE_SIZE + (long) p;  // Kernel mode stack pointer.
94         p->tss.ss0 = 0x10;              // Stack segment selector (same as data segment).
95         p->tss.eip = eip;
96         p->tss.eflags = eflags;
97         p->tss.eax = 0;                 // This is the reason why child return 0.
98         p->tss.ecx = ecx;
99         p->tss.edx = edx;
```

```
100         p->tss.ebx = ebx;
   // The new task duplicates the stack contents of the father completely, thus, we require
   // that the stack of the task 0 to be clean for the sake of this action.
101         p->tss.esp = esp;
102         p->tss.ebp = ebp;
103         p->tss.esi = esi;
104         p->tss.edi = edi;
105         p->tss.es = es & 0xffff;      // Segment register has only 16 bits.
106         p->tss.cs = cs & 0xffff;
107         p->tss.ss = ss & 0xffff;
108         p->tss.ds = ds & 0xffff;
109         p->tss.fs = fs & 0xffff;
110         p->tss.gs = gs & 0xffff;
   // Install the selector of the LDT descriptor. The LDT descriptor is in the GDT.
111         p->tss.ldt = _LDT(nr);
112         p->tss.trace_bitmap = 0x80000000;
   // If the current process has used coprocessor, then its context should be saved into
   // tss.i387 structure. The instruction 'clts' is used to clear the task switched flag TS
   // in the control register CR0. The TS flag helps to determine when the context of the
   // coprocessor does not match that of the task being executed by the CPU. Whenever a task
   // switched, the CPU sets the flag TS. If the TS set, then each 'ESC' instructions causes
   // exception 7 and is captured. The 'WAIT' instruction also causes exception 7 if both TS
   // and MP are set. Thus, when a task switch is happened after one instruction 'ESC', the
   // context of the coprocessor may need to be saved before executing a new 'ESC' instruction.
   // The exception handling procedure should save the context of the coprocessor and reset TS.
   // 'fnsave' instruction is used to save the context of the coprocessor into the memory area
   // (here is the tss.i387) specified by the target operand.
113         if (last_task_used_math == current)
114             __asm__("clts ; fnsave %0"::"m" (p->tss.i387));
   // Setup base addresses and limits in code and data segment descriptors for the new task.
   // Duplicate the page tables from that of the current process. If an error occurred, then
   // clear the task table slot for the new task, and free the page occupied by the new task
   // structure.
115         if (copy_mem(nr,p)) {            // Error occurred if the return value is non-zero.
116             task[nr] = NULL;
117             free_page((long) p);
118             return -EAGAIN;
119         }
   // If there is files opened in the fahter process, then increasing the counts by 1 for each
   // openning files respectively.
120         for (i=0; i<NR_OPEN;i++)
121             if (f=p->filp[i])
122                 f->f_count++;
   // Increasing current (father) process's counts of the pwd, root and executable by 1 each.
123         if (current->pwd)
124             current->pwd->i_count++;
125         if (current->root)
126             current->root->i_count++;
127         if (current->executable)
128             current->executable->i_count++;
   // Install the TSS and LDT segment descriptors in GDT for the new task. The limits of the
   // two segments are set to be 104 bytes equally. Refer to L52-66 in include/asm/system.h.
   // Notice that the task register TR will be loaded by the CPU automatically.
```

```
129         set_tss_desc(gdt+(nr<<1)+FIRST_TSS_ENTRY,&(p->tss));
130         set_ldt_desc(gdt+(nr<<1)+FIRST_LDT_ENTRY,&(p->ldt));
131         p->state = TASK_RUNNING;            /* do this last, just in case */
132         return last_pid;                    // Return the new task pid to the father.
133 }
134
    // Generating a non existing pid and find an empty slot (task number) in the task table
    // for the new task. the new pid is stored in a global variable 'last_pid'.
135 int find_empty_process(void)
136 {
137         int i;
138
139         repeat:
    // If the value of 'last_pid' is out of the range of a positive integer after increased
    // by 1, then recycling the pid number from 1 again. Afterwards, check the new last_pid
    // to see if it has been used by any tasks in the system, otherwise repeat the opeartion
    // above to get a new one.
140                 if ((++last_pid)<0) last_pid=1;
141                 for(i=0 ; i<NR_TASKS ; i++)
142                         if (task[i] && task[i]->pid == last_pid) goto repeat;
    // Scan the task table to get and return an empty slot index (task number).
143         for(i=1 ; i<NR_TASKS ; i++)            // Slot 0 is excluded.
144                 if (!task[i])
145                         return i;
    // If all the 64 slots in the task table are all occupied, return with an error code.
146         return -EAGAIN;
147 }
148
```

## 5.11.3  Information

### 5.11.3.1  Task State Segment (TSS)[1]

To provide efficient, protected multitasking, the 80386 employs several special data structures. It does not, however, use special instructions to control multitasking; instead, it interprets ordinary control-transfer instructions differently when they refer to the special data structures. The registers and data structures that support multitasking are:

- Task state segment
- Task state segment descriptor
- Task register
- Task gate descriptor

With these structures the 80386 can rapidly switch execution from one task to another, saving the context of the original task so that the task can be restarted later. In addition to the simple task switch, the 80386 offers two other task-management features:

1. Interrupts and exceptions can cause task switches (if needed in the system design). The processor not only switches automatically to the task that handles the interrupt or exception, but it automatically switches back to the interrupted task when the

---

[1] Copied and modified from Intel 80386 programming reference manual.

interrupt or exception has been serviced. Interrupt tasks may interrupt lower-priority interrupt tasks to any depth.

2. With each switch to another task, the 80386 can also switch to another LDT and to another page directory. Thus each task can have a different logical-to-linear mapping and a different linear-to-physical mapping. This is yet another protection feature, because tasks can be isolated and prevented from interfering with one another.

All the information the processor needs in order to manage a task is stored in a special type of segment, a task state segment (TSS). Figure 5-12 shows the format of a TSS for executing 80386 tasks.

| 31 23 | 15 7 0 | |
|---|---|---|
| I/O MAP BASE | 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | 64 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | LDT selector | 60 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | GS | 5C |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | FS | 58 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | DS | 54 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | SS | 50 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | CS | 4C |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | ES | 48 |
| EDI | | 44 |
| ESI | | 40 |
| EBP | | 3C |
| ESP | | 38 |
| EBX | | 34 |
| EDX | | 30 |
| ECX | | 2C |
| EAX | | 28 |
| EFLAGS | | 24 |
| INSTRUCTION POINTER (EIP) | | 20 |
| PAGE DIRECTORY REGISTER CR3 (PDBR) | | 1C |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | SS2 | 18 |
| ESP2 | | 14 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | SS1 | 10 |
| ESP1 | | 0C |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | SS0 | 08 |
| ESP0 | | 04 |
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | BACK LINK TO PREVIOUS TSS | 00 |

Figure 5-12 Intel 80386 32-bit Task State Segment

The fields of a TSS belong to two classes:

1. A dynamic set that the processor updates with each switch from the task. This set includes the fields that store:

   ◆The general registers (EAX, ECX, EDX, EBX, ESP, EBP, ESI, EDI).

* The segment registers (ES, CS, SS, DS, FS, GS).
* The flags register (EFLAGS).
* The instruction pointer (EIP).
* The selector of the TSS of the previously executing task (updated only when a return is expected).

2. A static set that the processor reads but does not change. This set includes the fields that store:
* The selector of the task's LDT.
* The register (PDBR) that contains the base address of the task's page directory (read only when paging is enabled).
* Pointers to the stacks for privilege levels 0-2.
* The T-bit (debug trap bit) which causes the processor to raise a debug exception when a task switch occurs . (Refer to Chapter 12 for more information on debugging.)
* The I/O map base (refer to Chapter 8 for more information on the use of the I/O map).

Task state segments may reside anywhere in the linear space. The only case that requires caution is when the TSS spans a page boundary and the higher-addressed page is not present. In this case, the processor raises an exception if it encounters the not-present page while reading the TSS during a task switch. Such an exception can be avoided by either of two strategies:

1. By allocating the TSS so that it does not cross a page boundary.
2. By ensuring that both pages are either both present or both not-present at the time of a task switch. If both pages are not-present, then the page-fault handler must make both pages present before restarting the instruction that caused the task switch.

The state segment, like all other segments, is defined by a descriptor. The task register (TR) identifies the currently executing task by pointing to the TSS. The instructions LTR and STR are used to modify and read the visible portion of the task register.

The I/O instructions that directly refer to addresses in the processor's I/O space are IN, INS, OUT, OUTS. The 80386 has the ability to selectively trap references to specific I/O addresses. The structure that enables selective trapping is the I/O Permission Bit Map in the TSS segment. The I/O permission map is a bit vector. The size of the map and its location in the TSS segment are variable. The processor locates the I/O permission map by means of the I/O map base field in the fixed portion of the TSS. The I/O map base field is 16 bits wide and contains the offset of the beginning of the I/O permission map. The upper limit of the I/O permission map is the same as the limit of the TSS segment.

In protected mode, when it encounters an I/O instruction (IN, INS, OUT, or OUTS), the processor first checks whether CPL <= IOPL (in EFLAGS) . If this condition is true, the I/O operation may proceed. If not true, the processor checks the I/O permission map.

Each bit in the map corresponds to an I/O port byte address; for example, the bit for port 41 is found at I/O map base + 5, bit offset 1. The processor tests all the bits that correspond to the I/O addresses spanned by an I/O operation; for example, a doubleword operation tests four bits corresponding to four adjacent byte addresses. If any tested bit is set, the processor signals a general protection exception. If all the tested bits are zero, the I/O operation may proceed.

It is not necessary for the I/O permission map to represent all the I/O addresses. I/O

addresses not spanned by the map are treated as if they had one bits in the map. For example, if TSS limit is equal to I/O map base + 31, the first 256 I/O ports are mapped; I/O operations on any port greater than 255 cause an exception. If I/O map base is greater than or equal to TSS limit, the TSS segment has no I/O permission map, and all I/O instructions in the 80386 program cause exceptions when CPL > IOPL. In the case of Linux 0.11 kernel, The I/O map base is set to be 0x8000, greater than that of the TSS limit (104 bytes). So the TSS has no I/O permission map.

Because the I/O permission map is in the TSS segment, different tasks can have different maps. Thus, the operating system can allocate ports to a task by changing the I/O permission map in the task's TSS.

In Linux 0.11, SS0:ESP0 in the TSS is used to store the kernel mode stack pointer of the task. Whenever a task enter into the kernel, the initial location of the kernel mode stack pointer is not changed, SS0:ESP0 will remain points to the top boundary of the memory page used by the task data structure. SS1:ESP1 and SS2:ESP2 are used as stack pointers for code running under level 1 and level 2. The Linux kernel does not use these two levels. The user mode stack pointer is stored in registers SS:ESP.

## 5.12  sys.c

### 5.12.1  Function

sys.c program includes implementations of many system calls. A few of the system calls have not been implementated in kernel 0.11 yet, and only renturn an error code -ENOSYS. A simple descriptions of all the system calls can refer to file include/linux/sys.h.

There are various kind of identifiers (IDs) in the program, including process ID (pid), process group ID (pgrp or pgid), user ID (uid), user group ID (gid), real user ID (ruid), real group ID (rgid), effective user ID (euid), effective group ID (egid) and session ID. Understanding them is important to read the code.

A user has a user ID (uid) and a group ID (gid), which are set in the passwd file for the user, and usually called real user ID (ruid) and real group ID (rgid). For each file in a file system, the I-node contains the uid and gid, which specify the owner of the file and what user group the file belongs to, and are used for the file access permission. In addition, a process contains three kinds of IDs in its task structure as shown in Table 5-4.

Table 5-4 User ID and group ID related to process.

| Type | User ID | Group ID |
|------|---------|----------|
| Process's | uid – User ID, specifies the owner of the process. | gid – Group ID, specifies the user group of the process. |
| Effective | euid – Effective user ID, indicates the file access permission. | egid – Effective goup Id, indicates the file access permission. |
| Saved | suid – Saved user ID. When the set-user-ID flag of a file is set, the suid is equal to the uid of the file, otherwise the suid | sgid – Saved group ID. When the set-group-ID flag of a file is set, the sgid is equal to the gid of the file, otherwise the sgid equals the |

equals the euid of the process.          egid of the process.

The uid and gid are that of the owner of the process, and they are actually the real user ID (ruid) and real group ID (rgid) of the process. User root (supervisor) may use functions set_uid() and set_gid() to modify them. The effective user ID (euid) and effective group ID (egid) are used to check the file access permission.

Saved user ID (suid) and saved group ID (sgid) are used to access file, which has the set-user-ID or set-group-ID flag set. When running a program, the process's euid is usually equal to the uid, and process's egid equals to gid. Thus, the process can only access the files which owns by the same euid and egid, or other files permitted. But, if a file's set-user-ID flag is set, then the euid of the process will be set to equal to the uid of the file owner. Therefore the process are permitted to access the restricted file and the uid of the file owner is saved in suid. In the same reason, the set-group-ID flag has the alike property and be operated in the similar method.

If, for example, the owner of a program file is user root (value is 0) and its set-user-ID flag is set, then, when the program is executed in a process, the process's effective user ID (euid) will be set to be root's ID (0). Thus, the process owns the root authority. A practical example is the line command passwd owned by user root. This command is a set-user-ID program and, therefore, allows normal users modifying their login passwords. Since the command needs to modify the password in the file /etc/passwd, which is also owned by the root, thus the passwd command executed by a normal user needs set-user-ID flag to be set.

In addition, A process has its own IDs, the process ID (pid), the process group ID (pgrp or pgid) and the session ID (session). These three IDs are used to indicate the relationship between processes and have nothing to do with user ID and user group ID.

## 5.12.2  Comments

Program 5-10 linux/kernel/sys.c

```
1  /*
2   *  linux/kernel/sys.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  #include <errno.h>
8
9  #include <linux/sched.h>
10 #include <linux/tty.h>
11 #include <linux/kernel.h>
12 #include <asm/segment.h>
13 #include <sys/times.h>
14 #include <sys/utsname.h>
15
   // Return current date and time.
   // Unimplemented system call and always return -1 and set 'errno' to ENOSYS.
16 int sys_ftime()
17 {
```

```
18          return -ENOSYS;                         // Not implemented yet.
19 }
20
   // Unimplemented system call. Obsoleted.
21 int sys_break()
22 {
23          return -ENOSYS;
24 }
25
   // The ptrace system call provides a means by which a parent process may observe and control
   // the execution of another process, and examine and change its core image and registers.
   // It is primarily used to implement breakpoint debugging and system call tracing.
26 int sys_ptrace()
27 {
28          return -ENOSYS;
29 }
30
   // Change and print terminal line settings.
31 int sys_stty()
32 {
33          return -ENOSYS;
34 }
35
   // Get terminal line settings. Obsoleted.
36 int sys_gtty()
37 {
38          return -ENOSYS;
39 }
40
   // Change the name or location of a file.
41 int sys_rename()
42 {
43          return -ENOSYS;
44 }
45
   // Unimplemented system call. Obsoleted.
46 int sys_prof()
47 {
48          return -ENOSYS;
49 }
50
   // Set real and/or effective group ID of the current process.
   // Unprivileged users may only interchange the real group ID and the effective group ID,
   // and may only set the effective group ID to the saved group ID.
51 int sys_setregid(int rgid, int egid)
52 {
53          if (rgid>0) {
54                  if ((current->gid == rgid) ||
55                     suser())
56                      current->gid = rgid;
57                  else
58                      return(-EPERM);
59          }
```

```
60              if (egid>0) {
61                      if ((current->gid == egid) ||
62                          (current->egid == egid) ||
63                          (current->sgid == egid) ||
64                          suser())
65                              current->egid = egid;
66                      else
67                              return(-EPERM);
68              }
69              return 0;
70      }
71
        // Set group identity (gid). Set the effective group ID of the current process. If the
        // the caller is the superuser, the real and saved group ID's are also set.
72      int sys_setgid(int gid)
73      {
74              return(sys_setregid(gid, gid));
75      }
76
        // Switch process accounting of or off.
77      int sys_acct()
78      {
79              return -ENOSYS;
80      }
81
        // Map any physical memory area into process's virtual space.
82      int sys_phys()
83      {
84              return -ENOSYS;
85      }
86
87      int sys_lock()
88      {
89              return -ENOSYS;
90      }
91
92      int sys_mpx()
93      {
94              return -ENOSYS;
95      }
96
97      int sys_ulimit()
98      {
99              return -ENOSYS;
100     }
101
        // Return the time since the Epoch (00:00:00 UTC, January 1, 1970), measured in seconds.
        // If tloc is non-NULL, the return value is also stored in the memory pointed to by it.
102     int sys_time(long * tloc)
103     {
104             int i;
105
106             i = CURRENT_TIME;
```

```
107          if (tloc) {
108                  verify_area(tloc,4);              // Verify the memory space. (4 bytes)
109                  put_fs_long(i,(unsigned long *)tloc); // Stored in tloc in user space.
110          }
111          return i;
112 }
113
114 /*
115  * Unprivileged users may change the real user id to the effective uid
116  * or vice versa.
117  */
    // setreuid sets real and effective user IDs of the current process. Unprivileged users
    // may only set the real user ID to the real user ID or the effective user ID, and may
    // only set the effective user ID to the real user ID, the effective user ID or the saved
    // user ID.
    // If the real user ID is set or the effective user ID is set to a value not equal to the
    // previous real user ID, the saved user ID will be set to the new effective user ID.
    // Completely  analogously, setregid sets real and effective group ID's of the current
    // process, and all of the above holds with "group" instead of "user".
118 int sys_setreuid(int ruid, int euid)
119 {
120          int old_ruid = current->uid;
121
122          if (ruid>0) {
123                  if ((current->euid==ruid) ||
124                      (old_ruid == ruid) ||
125                      suser())
126                          current->uid = ruid;
127                  else
128                          return(-EPERM);
129          }
130          if (euid>0) {
131                  if ((old_ruid == euid) ||
132                      (current->euid == euid) ||
133                      suser())
134                          current->euid = euid;
135                  else {
136                          current->uid = old_ruid;
137                          return(-EPERM);
138                  }
139          }
140          return 0;
141 }
142
    // setuid() sets the effective user ID of the current process. Unprivileged user may use
    // setuid() to set the effective user ID to the saved user ID, or real user ID. If the
    // caller is superuser, then the real, effective and saved user ID are all set to specified
    // uid.
143 int sys_setuid(int uid)
144 {
145          return(sys_setreuid(uid, uid));
146 }
147
```

```
      // stime sets the system's idea of the time and date. Time, pointed to by tptr, is
      // measured in seconds from 00:00:00 GMT January 1, 1970. stime() may only be executed by
      // the super user.
148 int sys_stime(long * tptr)
149 {
150        if (!suser())
151                return -EPERM;
152        startup_time = get_fs_long((unsigned long *)tptr) - jiffies/HZ;
153        return 0;
154 }
155
      // Get process times. Store the current process times user time (utime), sys time (stime),
      // child user time (cutime) and child sys time (cstime) in the struct tms that tbuf points to.
156 int sys_times(struct tms * tbuf)
157 {
158        if (tbuf) {
159                verify_area(tbuf,sizeof *tbuf);
160                put_fs_long(current->utime, (unsigned long *)&tbuf->tms_utime);
161                put_fs_long(current->stime, (unsigned long *)&tbuf->tms_stime);
162                put_fs_long(current->cutime, (unsigned long *)&tbuf->tms_cutime);
163                put_fs_long(current->cstime, (unsigned long *)&tbuf->tms_cstime);
164        }
165        return jiffies;
166 }
167
      // brk sets the end of the data segment to the value specified by 'end_data_seg', when that
      // value is reasonable, the system does have enough memory and the process does not exceed
      // its max data size. The value must be greater than the end of code and less than the 16KB
      // from the end the stack. The return value is the new end value of the data seg.
      // This system call is usually packaged into a function in the library libc and the return
      // value is different.
168 int sys_brk(unsigned long end_data_seg)
169 {
      // If the argument 'end_data_seg' greater than the end of code and less than the start
      // point of the stack - 16KB, then set the end of data segment in the current process.
170        if (end_data_seg >= current->end_code &&
171            end_data_seg < current->start_stack - 16384)
172                current->brk = end_data_seg;
173        return current->brk;
174 }
175
176 /*
177  * This needs some heave checking ...
178  * I just haven't get the stomach for it. I also don't fully
179  * understand sessions/pgrp etc. Let somebody who does explain it.
180  */
      // setpgid sets the process group ID of the process specified by pid to pgid. If pid is
      // zero, the process ID of the current process is used. If pgid is zero, the process ID
      // of the process specified by pid is used. If setpgid is used to move a process from one
      // process group to another, both process groups must be part of the same session. In this
      // case, the pgid specifies an existing process group to be joined and the session ID of
      // that group must match the session ID of the joining process (L193).
      // getpgid returns the process group ID of the process specified by pid. If pid is zero,
```

```
    // the process ID of the current process is used.
181 int sys_setpgid(int pid, int pgid)
182 {
183         int i;
184
185         if (!pid)                           // If pid=0, uses pid of the current process.
186                 pid = current->pid;
187         if (!pgid)                          // If pgid=0, uses the current's pid as pgid.
188                 pgid = current->pid;        // [Here is iffer from the description of POSIX]
189         for (i=0 ; i<NR_TASKS ; i++)        // Scan task Table, search the process pid.
190                 if (task[i] && task[i]->pid==pid) {
191                         if (task[i]->leader)        // Already leader?
192                                 return -EPERM;
193                         if (task[i]->session != current->session)
194                                 return -EPERM;
195                         task[i]->pgrp = pgid;       // Set task's pgrp.
196                         return 0;
197                 }
198         return -ESRCH;
199 }
200
    // Return the process group ID of the current process. Same as getpgid(0).
201 int sys_getpgrp(void)
202 {
203         return current->pgrp;
204 }
205
    // setsid() creates a new session if the calling process is not a process group leader.
    // The calling process is the leader of the new session, the process group leader of the
    // new process group, and has no controlling tty. The process group ID and session ID of
    // the calling process are set to the PID of the calling process. The calling process will
    // be the only process in this new process group and in this new session.
    // setsid -- SET Session ID。
206 int sys_setsid(void)
207 {
    // If the current process is not a process group leader and is not owned by a superuser,
    // then return error code.
208         if (current->leader && !suser())
209                 return -EPERM;
    // Set the current process to be the session leader, and its session = pgrp = pid.
210         current->leader = 1;
211         current->session = current->pgrp = current->pid;
212         current->tty = -1;                      // The process has no controlling terminal.
213         return current->pgrp;                   // Return session ID.
214 }
215
    // uname() returns system information in the structure pointed to by name.
    // The utsname struct is defined in sys/utsname.h
216 int sys_uname(struct utsname * name)
217 {
    // Here gives the information in the structure. This coding style will surely be changed.
218         static struct utsname thisname = {
219                 "linux .0","nodename","release ","version ","machine "
```

```
220            };
221            int i;
222
223            if (!name) return -ERROR;          // Error if the buffer pointer is NULL.
224            verify_area(name,sizeof *name);   // Verify the memory area for saving info.
225            for(i=0;i<sizeof *name;i++)       // Copy the info to user data space.
226                    put_fs_byte(((char *) &thisname)[i],i+(char *) name);
227            return 0;
228 }
229
    // Set the file creation mask to be mask & 0777, and return the orignal mask.
230 int sys_umask(int mask)
231 {
232            int old = current->umask;
233
234            current->umask = mask & 0777;
235            return (old);
236 }
237
```

# 5.13   vsprintf.c

## 5.13.1   Function

function vsprintf() is used to produce formatted output to the character string, and is one function in the printf() family. Since it is a function in the standard library now and with almost no relationship to the remain kernel code, you can refer to the library manual to understand the usage of it and might skip it now.

## 5.13.2   Comments

Program 5-11 linux/kernel/vsprintf.c

```
 1 /*
 2  *  linux/kernel/vsprintf.c
 3  *
 4  *  (C) 1991  Linus Torvalds
 5  */
 6
 7 /* vsprintf.c -- Lars Wirzenius & Linus Torvalds. */
 8 /*
 9  * Wirzenius wrote this portably, Torvalds fucked it up :-)
10  */
    // Lars Wirzenius is one of the best friends of Linus at Helsinki university, and stayed in
    // the same office once. When developing the Linux operation system in 1991, Linus is not
    // very skillful with C language and know little about the variable-length argument facilities.
    // Thus, Lars Wirzenius helped programmed this code. This program is much robust and only
    // one bug founded 3 year later (Line 130 bellow). His homepage: http://liw.iki.fi/liw/
11
12 #include <stdarg.h>
13 #include <string.h>
```

```
14
15  /* we use this so that we can do without the ctype library */
16  #define is_digit(c)      ((c) >= '' && (c) <= '9')
17

    // Convert a string to value. The argument is a pointer to pointer to a string.
    // The returns is the converted value. In addition, the pointer steps forward.
18  static int skip_atoi(const char **s)
19  {
20          int i=0;
21
22          while (is_digit(**s))
23                  i = i*10 + *((*s)++) - '';
24          return i;
25  }
26

    // Defines the symbols for the converting type.
27  #define ZEROPAD 1             /* pad with zero */
28  #define SIGN    2             /* unsigned/signed long */
29  #define PLUS    4             /* show plus */
30  #define SPACE   8             /* space if plus */
31  #define LEFT    16            /* left justified */
32  #define SPECIAL 32            /* 0x */
33  #define SMALL   64            /* use 'abcdef' instead of 'ABCDEF' */
34

    // Divide operation.
    // Input: 'n' is a dividend, 'base' is a divisor.
    // The result 'n' is quotient. The return value is the remainder.
35  #define do_div(n,base) ({ \
36  int __res; \
37  __asm__("divl %4":"=a" (n),"=d" (__res):"" (n),"1" (0),"r" (base)); \
38  __res; })
39

    // Convert a integer to a string with specified base system.
    // Arguments: num - Integer value; base - system base; size - length of string;
    //            precision - digital length; type - type option (Refer to L27-33).
    // Output: A pointer to the string str.
40  static char * number(char * str, int num, int base, int size, int precision
41          ,int type)
42  {
43          char c,sign,tmp[36];
44          const char *digits="0123456789ABCDEFGHIJKLMNOPQRSTUVWXYZ";
45          int i;
46

    // If the type indicates lowercase, then defines the lowercase character set.
    // If the type indicates left justified, then clear the flag 'ZEROPAD' in type.
    // If the carry base is out of the range 2 and 36 then exit.
47          if (type&SMALL) digits="0123456789abcdefghijklmnopqrstuvwxyz";
48          if (type&LEFT) type &= ~ZEROPAD;
49          if (base<2 || base>36)
50                  return 0;
    // If the type indicates the need of zero padded, the set variable c='0' (that is ''),
    // otherwise set c equals space character (c=' ').
    // If the type indicates it's a signed value and less than 0, then set sign character and
```

```
      // get the absolute value of the input integer. Otherwise if the type indicates a sign of
      // plus, then set a plus character to sign, and set sign to space character if indicated
      // or set to '0'.
51            c = (type & ZEROPAD) ? '' : ' ';
52            if (type&SIGN && num<0) {
53                    sign='-';
54                    num = -num;
55            } else
56                    sign=(type&PLUS) ? '+' : ((type&SPACE) ? ' ' : 0);
      // If it's a signed integer, then decrease its szie with 1. If it's a special convertion,
      // then minus 2 once more for hex (used for the prefix '0x') or minus 1 once more for the
      // octals ( for the prefix '0').
57            if (sign) size--;
58            if (type&SPECIAL)
59                    if (base==16) size -= 2;
60                    else if (base==8) size--;
      // If the integer is 0, then set temporary string array temp to '0', otherwise do the
      // convertion operation using specified base system.
61            i=0;
62            if (num==0)
63                    tmp[i++]='';
64            else while (num!=0)
65                    tmp[i++]=digits[do_div(num,base)];
      // If the number of the digital charactors great than the value of precision, then extends
      // the length of precision to match the number digital charactors.
66            if (i>precision) precision=i;
67            size -= precision;

      // Here begins the real convertion and store the result into string 'str' temporarily.
      // If the 'type' does not contain flag 'ZEROPAD' and 'LEFT', then we first fill the space
      // charactor into 'str' with number of remain size. If it has sign symbol, then fill the
      // sign character.
68            if (!(type&(ZEROPAD+LEFT)))
69                    while(size-->0)
70                            *str++ = ' ';
71            if (sign)
72                    *str++ = sign;
      // If the type indicates that it is a special conversion, the put a '0' at the beginning
      // of str for the ocatals or put '0x' for the hex.
73            if (type&SPECIAL)
74                    if (base==8)
75                            *str++ = '';
76                    else if (base==16) {
77                            *str++ = '';
78                            *str++ = digits[33];    // 'X'或'x'
79                    }
      // If there is no 'LEFT' flag in the 'type', then fill with paddle charactors ('0' or
      // space character).
80            if (!(type&LEFT))
81                    while(size-->0)
82                            *str++ = c;
      // At this moment, the variable 'i' contains the number of digitals for 'num'. If the
      // number is less than the size of precision, then fills 'str' with precison-i number
```

```
       // of character '0'.
83         while(i<precision--)
84                 *str++ = '';
       // Fill the 'str' with converted digitals, i totally.
85         while(i-->0)
86                 *str++ = tmp[i];
       // If the size is still great than 0, means there is a 'LEFT' flag in the type. So we
       // put space characters in the remain size of the 'str'.
87         while(size-->0)
88                 *str++ = ' ';
89         return str;
90 }
91
       // Output the formatted string.
       // The arguments fmt is the formatting string, args is variable arguments, buf is the
       // buffer for formatted string.
92 int vsprintf(char *buf, const char *fmt, va_list args)
93 {
94         int len;
95         int i;
96         char * str;               // Used for temporarily store the string during conversion.
97         char *s;
98         int *ip;
99
100        int flags;               /* flags to number() */
101
102        int field_width;         /* width of output field */
103        int precision;           /* min. # of digits for integers; max
104                                    number of chars for from string */
105        int qualifier;           /* 'h', 'l', or 'L' for integer fields */
106
       // First we points to buffer 'buf', then scan the formating string 'fmt', obtaining the
       // various converting flags.
107        for (str=buf ; *fmt ; ++fmt) {
108                if (*fmt != '%') {
109                        *str++ = *fmt;
110                        continue;
111                }
112
113                /* process flags */
114                flags = 0;
115                repeat:
116                        ++fmt;           /* this also skips first '%' */
117                        switch (*fmt) {
118                                case '-': flags |= LEFT; goto repeat;
119                                case '+': flags |= PLUS; goto repeat;
120                                case ' ': flags |= SPACE; goto repeat;
121                                case '#': flags |= SPECIAL; goto repeat;
122                                case '': flags |= ZEROPAD; goto repeat;
123                        }
124
       // Get width field of the current argument, and put into the variable 'field_width'.
       // If the width field contains digitals, then set it to be the width value. If it is
```

```
        // character '*', it means the next argument specifies the witdh. Thus, we invoke the macro
        // va_arg to get the width value. If, at this moment, the width value is less than 0, then
        // it means there is a flag field (Left justified). Thereby we should append this flag to
        // the flag variable 'flag' and set field width with its absolute value.
125                     /* get field width */
126                     field_width = -1;
127                     if (is_digit(*fmt))
128                             field_width = skip_atoi(&fmt);
129                     else if (*fmt == '*') {
130                             /* it's the next argument */    // A bug. Should insert '++fmt;'
131                             field_width = va_arg(args, int);
132                             if (field_width < 0) {
133                                     field_width = -field_width;
134                                     flags |= LEFT;
135                             }
136                     }
137

        // Get the precision field. The precision field begins with '.'. The handling way is
        // similar to that of the width field. If the precision field contains digitals, then
        // set it to be the precision value. If it is a character '*', means the next argument
        // specifies the precision. So we invoke the macro va_arg to get the precision. If the
        // value is less than 0, we get its absolute as its precision value.
138                     /* get the precision */
139                     precision = -1;
140                     if (*fmt == '.') {
141                             ++fmt;
142                             if (is_digit(*fmt))
143                                     precision = skip_atoi(&fmt);
144                             else if (*fmt == '*') {
145                                     /* it's the next argument */
146                                     precision = va_arg(args, int);
147                             }
148                             if (precision < 0)
149                                     precision = 0;
150                     }
151
        // The following code analyzes the conversion gualifier.
        // 'h', 'l', 'L' are qualifers for input length modifier.
152                     /* get the conversion qualifier */
153                     qualifier = -1;
154                     if (*fmt == 'h' || *fmt == 'l' || *fmt == 'L') {
155                             qualifier = *fmt;
156                             ++fmt;
157                     }
158
        // The following code analyzes the conversion specifier.
159                     switch (*fmt) {
        // If the specifier is 'c', it means the corresponding argument is a character. At this
        // moment, if the flag indicates with no left justified, then we should put number of
        // 'width value -1' space characters ahead of the field. If the width value is still great
        // than 0, means left justified, then we put number of 'width value - 1' space characters
        // after the field.
160                     case 'c':
```

```
161                         if (!(flags & LEFT))
162                                 while (--field_width > 0)
163                                         *str++ = ' ';
164                         *str++ = (unsigned char) va_arg(args, int);
165                         while (--field_width > 0)
166                                 *str++ = ' ';
167                         break;
168
```

// If the specifier is character 's', it means the corresponding argument is a string. First,
// we obtain the length of the string. If it great than the precision field value, then we
// extend the precision value to be equal to the string length. If, at thim moment, the flag
// does not indicate left justified, then we should put number of 'field width - string len'
// of space characters ahead of the field. If the width value is still great than 0, means it
// needs left justified, then we put number of 'field width - string len ' space characters.

```
169                 case 's':
170                         s = va_arg(args, char *);
171                         len = strlen(s);
172                         if (precision < 0)
173                                 precision = len;
174                         else if (len > precision)
175                                 len = precision;
176
177                         if (!(flags & LEFT))
178                                 while (len < field_width--)
179                                         *str++ = ' ';
180                         for (i = 0; i < len; ++i)
181                                 *str++ = *s++;
182                         while (len < field_width--)
183                                 *str++ = ' ';
184                         break;
185
```

// If the specifier is 'o', it means the corresponding argument needs convert to string
// in octals.

```
186                 case 'o':
187                         str = number(str, va_arg(args, unsigned long), 8,
188                                 field_width, precision, flags);
189                         break;
190
```

// If the specifier is 'p', it means the corresponding argument is a pointer. In this case,
// if the argument has no width field, the default width is 8, and need to append 0.

```
191                 case 'p':
192                         if (field_width == -1) {
193                                 field_width = 8;
194                                 flags |= ZEROPAD;
195                         }
196                         str = number(str,
197                                 (unsigned long) va_arg(args, void *), 16,
198                                 field_width, precision, flags);
199                         break;
200
```

// If the specifier is 'x' or 'X', it means the argument needs print in hex. 'x' indicates
// print in lowercase.

```
201                 case 'x':
```

```
202                          flags |= SMALL;
203              case 'X':
204                      str = number(str, va_arg(args, unsigned long), 16,
205                              field_width, precision, flags);
206                  break;
207
```
// If the specifier is 'd', 'i' or 'u', they mean the argument is a integer. As 'd' and 'i'
// represents signed integer, a sign character need to appended.
```
208              case 'd':
209              case 'i':
210                      flags |= SIGN;
211              case 'u':
212                      str = number(str, va_arg(args, unsigned long), 10,
213                              field_width, precision, flags);
214                  break;
215
```
// If the specifer is 'n', it means we need to store the number of converted characters into
// the location pointed by the corresponding argument. Firstly, we use va_arg() to get the
// argument pointer, then save the number of converted characters into the location.
```
216              case 'n':
217                      ip = va_arg(args, int *);
218                      *ip = (str - buf);
219                  break;
220
```
// If the specifier is not character '%', it means there is an error in the formating string,
// then we put a '%' into the output string instead. If the specifier contains other more
// characters, they are put into the output string too.
```
221              default:
222                      if (*fmt != '%')
223                              *str++ = '%';
224                      if (*fmt)
225                              *str++ = *fmt;
226                      else
227                              --fmt;
228                  break;
229          }
230      }
231      *str = '\0';          // Append a NULL at the end of the converted string.
232      return str-buf;        // Return the length of the converted string.
233 }
234
```

# 5.14  printk.c

## 5.14.1  Function

printk() used in kernel code and has the same operation as printf() implemented in the
standard library. The reason to re-program it is that the kernel code is forbidden to modify
the segment register fs, which is dedicated to use in user mode. printk() uses function svprintf()

to format the arguments and then, after saved the register fs, invokes tty_write() to display messages on the console screen.

## 5.14.2 Commnets

<div align="center">Program 5-12 linux/kernel/printk.c</div>

```
1  /*
2   *  linux/kernel/printk.c
3   *
4   *  (C) 1991  Linus Torvalds
5   */
6
7  /*
8   * When in kernel-mode, we cannot use printf, as fs is liable to
9   * point to 'interesting' things. Make a printf with fs-saving, and
10  * all is well.
11  */
12 #include <stdarg.h>
13 #include <stddef.h>
14
15 #include <linux/kernel.h>
16
17 static char buf[1024];
18
19 extern int vsprintf(char * buf, const char * fmt, va_list args);
20
   // print function used in kernel code exclusively.
21 int printk(const char *fmt, ...)
22 {
23        va_list args;                    // va_list is really a character pointer.
24        int i;
25
   // Begin arguments processing. Output the formatted string into buffer. variable 'i' contains
   // the length of the formatted string.
26        va_start(args, fmt);
27        i=vsprintf(buf,fmt,args);
28        va_end(args);
   // Save the segment register fs, push three arguments for tty_write() and call it.
   // Afterwards restore the register fs and return the length of the string.
29        __asm__("push %%fs\n\t"          // Save fs
30                "push %%ds\n\t"
31                "pop %%fs\n\t"           // Set fs = ds
32                "pushl %0\n\t"           // Arg1: Length of string.
33                "pushl $_buf\n\t"        // Arg2: Address of the string buffer.
34                "pushl $0\n\t"           // Arg3: Channel No. (0).
35                "call _tty_write\n\t"    // (L290 in file kernel/chr_drv/tty_io.c)
36                "addl $8,%%esp\n\t"      // Discard 2 pushed args (buf, channel).
37                "popl %0\n\t"            // Pop the string length as return value.
38                "pop %%fs"               // Restore fs.
39                ::"r" (i):"ax","cx","dx");
40        return i;                        // The length of the string.
41 }
42
```

# 5.15 panic.c

## 5.15.1 Function

Function panic() is used to show the kernel error message and let the kernel code goes into deadly circle whenever kernel code goes wrong. This is a simple but efficient way to handle the kernel code bugs, and follow the prinicple of the UN*X systems: The simple, the best.

In the book 'Hitch hikers Guide to the Glaxy', written by Douglas Adams, one of the famous sentence is: Don't Panic!

## 5.15.2 Comments

程序 5-1 linux/kernel/panic.c

```
 1 /*
 2  *  linux/kernel/panic.c
 3  *
 4  *  (C) 1991  Linus Torvalds
 5  */
 6
 7 /*
 8  * This function is used through-out the kernel (includeinh mm and fs)
 9  * to indicate a major problem.
10  */
11 #include <linux/kernel.h>
12 #include <linux/sched.h>
13
14 void sys_sync(void);    /* it's really int */    /* (L44 in file fs/buffer.c) */
15
16 volatile void panic(const char * s)
17 {
18         printk("Kernel panic: %s\n\r",s);
19         if (current == task[0])
20                 printk("In swapper task - not syncing\n\r");
21         else
22                 sys_sync();
23         for(;;);
24 }
25
```

# 5.16 Summary

This chapter has disscused the 12 programs in the kernel subdirectory, which are the most important files in the kernel. The main contents include system calls, process schedule, signal processing, process forking and terminating.