

Week6 Report in Class (Fri56)

11911839 聂雨荷

Q1

代码中通过何种方式从S mode 进入U mode?

1. Create one process execute function `user_main()` in `init_main()` function we have discussed in the previous lab

- We create `proc_struct` for the new process, and add it to the process linked list

```
1 static int init_main(void *arg) {
2     size_t nr_free_pages_store = nr_free_pages();
3     //size_t kernel_allocated_store = kallocated();
4
5     int pid = kernel_thread(user_main, NULL, 0);
6     if (pid <= 0) {
7         panic("create user_main failed.\n");
8     }
9
10    ...
11 }
```

2. As the administrator, `init_main()`, will call `do_wait` to wait for the ending of the child process, and use `schedule()` to schedule `user_main()` during the wait

```
1 static int init_main(void *arg) {
2     ...
3
4     while (do_wait(0, NULL) == 0) {
5         schedule();
6     }
7
8     cprintf("all user-mode processes have quit.\n");
9     assert(initproc->cptr == NULL && initproc->yptr == NULL && initproc->optr == NULL);
10    assert(nr_process == 2);
11    assert(list_next(&proc_list) == &(initproc->list_link));
12    assert(list_prev(&proc_list) == &(initproc->list_link));
13
14    cprintf("init check memory pass.\n");
15    return 0;
16 }
```

3. As we have discussed in Lab5, `schedule()` function will create find the `user_main()` process and execute it. `user_main()` calls the `KERNEL_EXECVE()` function

```

1 // user_main - kernel thread used to exec a user program
2 static int user_main(void *arg) {
3     KERNEL_EXECVE(hello);
4     panic("user_main execve failed.\n");
5 }

```

```

1 #define __KERNEL_EXECVE(name, binary, size) ({ \
2     cprintf("kernel_execve: pid = %d, name = \"%s\".\n", \
3         current->pid, name); \
4     kernel_execve(name, binary, (size_t)(size)); \
5 })
6
7 #define KERNEL_EXECVE(x) ({ \
8     extern unsigned char _binary_obj__user_##x##_out_start[], \
9     _binary_obj__user_##x##_out_size[]; \
10    __KERNEL_EXECVE(#x, _binary_obj__user_##x##_out_start, \
11        _binary_obj__user_##x##_out_size); \
12 })

```

- This two function later will allows user_main() to execute the code hello.c

```

1 kern_execve("hello",
2 _binary_obj__user_hello_out_start, _binary_obj__hello_exit_out_size)

```

4. kern_execve() triggers ebreak exception

```

1 // kernel_execve - do SYS_exec syscall to exec a user program called by user_main
  kernel_thread
2 static int kernel_execve(const char *name, unsigned char *binary, size_t size) {
3     int64_t ret=0, len = strlen(name);
4     asm volatile(
5         "li a0, %1\n"
6         "lw a1, %2\n"
7         "lw a2, %3\n"
8         "lw a3, %4\n"
9         "lw a4, %5\n"
10        "li a7, 10\n"
11        "ebreak\n" // trigger exception
12        "sw a0, %0\n"
13        : "=m"(ret)
14        : "i"(SYS_exec), "m"(name), "m"(len), "m"(binary), "m"(size)
15        : "memory");
16    //cprintf("ret = %d\n", ret);
17    return ret;
18 }

```

5. This ebreak exception will be caught by trap() function. trap() function recognizes the cause is CAUSE_BREAKPOINT and enters syscall() to handle this problem

```

1 void exception_handler(struct trapframe *tf) {
2     int ret;
3     switch (tf->cause) {
4         ...
5         case CAUSE_BREAKPOINT:
6             cprintf("Breakpoint\n");

```

```

7         if(tf->gpr.a7 == 10){
8             //tf->epc += 4;
9             syscall();
10        }
11        break;
12        ...
13    }
14 }

```

6. `syscall()` calls the corresponding function defined by an array of syscalls, in this case, register `a0` has been set to be 4, which is `SYS_exec`, continuing execute `SYSexec()`

```

1 void
2 syscall(void) {
3     struct trapframe *tf = current->tf;
4     uint64_t arg[5];
5     int num = tf->gpr.a0;
6     if (num >= 0 && num < NUM_SYSCALLS) {
7         if (syscalls[num] != NULL) {
8             arg[0] = tf->gpr.a1;
9             arg[1] = tf->gpr.a2;
10            arg[2] = tf->gpr.a3;
11            arg[3] = tf->gpr.a4;
12            arg[4] = tf->gpr.a5;
13            tf->gpr.a0 = syscalls[num](arg);
14            return ;
15        }
16    }
17    panic("undefined syscall %d, pid = %d, name = %s.\n",
18        num, current->pid, current->name);
19 }
20
21
22 static int (*syscalls[])(uint64_t arg[]) = {
23     [SYS_exit]      sys_exit,
24     [SYS_fork]      sys_fork,
25     [SYS_wait]      sys_wait,
26     [SYS_exec]      sys_exec,
27     [SYS_yield]     sys_yield,
28     [SYS_kill]      sys_kill,
29     [SYS_getpid]    sys_getpid,
30     [SYS_putc]      sys_putc,
31 };

```

7. `sys_exec()` calls `do_execve()`

```

1 static int
2 sys_exec(uint64_t arg[]) {
3     const char *name = (const char *)arg[0];
4     size_t len = (size_t)arg[1];
5     unsigned char *binary = (unsigned char *)arg[2];
6     size_t size = (size_t)arg[3];
7     return do_execve(name, len, binary, size);
8 }

```

8. `do_execve()` allocates required memory for the new process

```

1 // do_execve - call exit_mmap(mm)&put_pgdir(mm) to reclaim memory space of current
  process
2 //          - call load_icode to setup new memory space according binary prog.
3 int
4 do_execve(const char *name, size_t len, unsigned char *binary, size_t size) {
5     struct mm_struct *mm = current->mm;
6     if (!user_mem_check(mm, (uintptr_t)name, len, 0)) {
7         return -E_INVALID;
8     }
9     if (len > PROC_NAME_LEN) {
10         len = PROC_NAME_LEN;
11     }
12
13     char local_name[PROC_NAME_LEN + 1];
14     memset(local_name, 0, sizeof(local_name));
15     memcpy(local_name, name, len);
16
17     if (mm != NULL) {
18         cputs("mm != NULL");
19         lcr3(boot_cr3);
20         if (mm_count_dec(mm) == 0) {
21             exit_mmap(mm);
22             put_pgdir(mm);
23             mm_destroy(mm);
24         }
25         current->mm = NULL;
26     }
27     int ret;
28     if ((ret = load_icode(binary, size)) != 0) { // prepare for the code
29         goto execve_exit;
30     }
31     set_proc_name(current, local_name);
32     return 0;
33
34 execve_exit:
35     do_exit(ret);
36     panic("already exit: %e.\n", ret);
37 }

```

9. load_icode() function sets the required resource. Specifically, it will set the tf->status to 0, which will wait for the setting of the operating system

```

1 /* load_icode - load the content of binary program(ELF format) as the new content of
  current process
2 * @binary: the memory addr of the content of binary program
3 * @size: the size of the content of binary program
4 */
5 static int
6 load_icode(unsigned char *binary, size_t size) {
7     ...
8     tf->status = sstatus & ~(SSTATUS_SPP | SSTATUS_SPIE);
9     ...
10    goto out;
11 }

```

10. after the trap finish loading the process, it will call the following code to return. This sret is a RISV command and sets CSRs[sstatus].SPP = 0, a User mode

```
1  # trapentry.S
2  __trapret:
3      RESTORE_ALL
4      # return from supervisor call
5      sret
6
7      .globl forkrets
```

Q2

代码中用户进程调用系统调用的过程是怎样的?

1. Initially, hello.c uses cprintf() to print the string

```
1  #include <stdio.h>
2  #include <ulib.h>
3
4  int
5  main(void) {
6      cprintf("Hello world!!.\n");
7      cprintf("I am process %d.\n", getpid());
8      cprintf("hello pass.\n");
9      return 0;
10 }
```

2. cprintf() encapsulates several functions. It finally calls the sys_putc() rather than call sbi_console_putchar(). This is because only S mode can use ecall to call the OpenSBI interface. We have to use system call.

```
1  /* *
2   * cprintf - formats a string and writes it to stdout
3   *
4   * The return value is the number of characters which would be
5   * written to stdout.
6   * */
7  int cprintf(const char *fmt, ...) {
8      va_list ap;
9
10     va_start(ap, fmt);
11     int cnt = vcprintf(fmt, ap); // call vcprintf
12     va_end(ap);
13
14     return cnt;
15 }
16
17 /* *
18 * vcprintf - format a string and writes it to stdout
19 *
20 * The return value is the number of characters which would be
21 * written to stdout.
22 *
23 * Call this function if you are already dealing with a va_list.
```

```

24  * Or you probably want cprintf() instead.
25  * */
26  int vcprintf(const char *fmt, va_list ap) {
27      int cnt = 0;
28      vprintfmt((void*)cputch, &cnt, fmt, ap); // call cputch
29      return cnt;
30  }
31
32  /* *
33   * cputch - writes a single character @c to stdout, and it will
34   * increace the value of counter pointed by @cnt.
35   * */
36  static void cputch(int c, int *cnt) {
37      sys_putc(c);
38      (*cnt) ++;
39  }

```

3. sys_putc() calls syscall() transform from U mode to S mode.

- sys_call() function makes ecall environment calls through inline assembly. This will generate a trap and enter S mode for exception handling

```

1  int sys_putc(int64_t c) {
2      return syscall(SYS_putc, c);
3  }
4
5  static inline int syscall(int64_t num, ...) {
6      va_list ap;
7      va_start(ap, num);
8      uint64_t a[MAX_ARGS];
9      int i, ret;
10     for (i = 0; i < MAX_ARGS; i++) {
11         a[i] = va_arg(ap, uint64_t);
12     }
13     va_end(ap);
14
15     asm volatile (
16         "ld a0, %1\n"
17         "ld a1, %2\n"
18         "ld a2, %3\n"
19         "ld a3, %4\n"
20         "ld a4, %5\n"
21         "ld a5, %6\n"
22         "ecall\n"
23         "sd a0, %0"
24         : "=m" (ret)
25         : "m"(num), "m"(a[0]), "m"(a[1]), "m"(a[2]), "m"(a[3]), "m"(a[4])
26         : "memory");
27     return ret;
28 }

```

4. Trap will be captured by trap.c and forward this system call to syscall() function defined in the kernel

```

1 void exception_handler(struct trapframe *tf) {
2     int ret;
3     switch (tf->cause) {
4         ...
5         case CAUSE_USER_ECALL:
6             //cprintf("Environment call from U-mode\n");
7             tf->epc += 4;
8             syscall();
9             break;
10        ...
11    }
12 }

```

5. `syscall()` takes the arguments in the register and forwards them to the function corresponding to the system call number for processing

```

1 void syscall(void) {
2     struct trapframe *tf = current->tf;
3     uint64_t arg[5];
4     int num = tf->gpr.a0;
5     if (num >= 0 && num < NUM_SYSCALLS) {
6         if (syscalls[num] != NULL) {
7             arg[0] = tf->gpr.a1;
8             arg[1] = tf->gpr.a2;
9             arg[2] = tf->gpr.a3;
10            arg[3] = tf->gpr.a4;
11            arg[4] = tf->gpr.a5;
12            tf->gpr.a0 = syscalls[num](arg);
13            return ;
14        }
15    }
16    panic("undefined syscall %d, pid = %d, name = %s.\n",
17          num, current->pid, current->name);
18 }

```

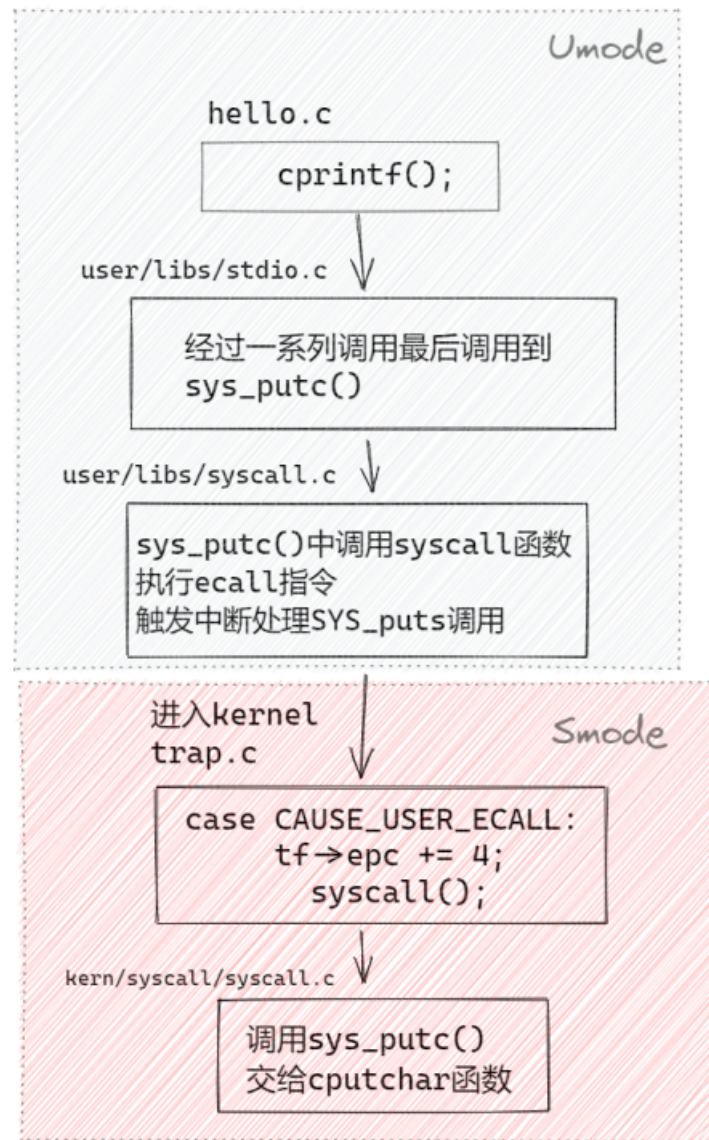
6. `syscall()` will trigger `sys_putc()`, as the user call requests system call to print the string

```

1 static int sys_putc(uint64_t arg[]) {
2     int c = (int)arg[0];
3     cputchar(c);
4     return 0;
5 }

```

7. We have discussed `cputchar()` function in Assignment2. It will trigger OpenSBI function. At this point, we have completed the system call.



Q3

代码中 hello.c 的 main() 执行结束后发生了什么，模式是否切换？

1. When the hello.c in the user mode finishes its execution, the process finishes its work. hello.c is actually called by umain.c, as the following code shows, umain will call the exit(ret)

```
1 void
2 umain(void) {
3     int ret = main();
4     exit(ret);
5 }
```

2. exit(ret) calls sys_exit, which is a system call


```

1 void
2 exit(int error_code) {
3     sys_exit(error_code);
4     cprintf("BUG: exit failed.\n");
5     while (1);
6 }

```

3. the rest of the procedures are like basic system call procedure we have discussed in the previous lab

- command `ecall` causes a trap, which will be caught and deliver to the kernel. At this time, **User mode is transformed into Kernel mode**

```

1 int
2 sys_exit(int64_t error_code) {
3     return syscall(SYS_exit, error_code);
4 }
5
6 static inline int
7 syscall(int64_t num, ...) {
8     va_list ap;
9     va_start(ap, num);
10    uint64_t a[MAX_ARGS];
11    int i, ret;
12    for (i = 0; i < MAX_ARGS; i++) {
13        a[i] = va_arg(ap, uint64_t);
14    }
15    va_end(ap);
16
17    asm volatile (
18        "ld a0, %1\n"
19        "ld a1, %2\n"
20        "ld a2, %3\n"
21        "ld a3, %4\n"
22        "ld a4, %5\n"
23        "ld a5, %6\n"
24        "ecall\n" // ecall, execute system call in the kernel mode
25        "sd a0, %0"
26        : "=m" (ret)
27        : "m"(num), "m"(a[0]), "m"(a[1]), "m"(a[2]), "m"(a[3]), "m"(a[4])
28        : "memory");
29    return ret;
30 }
31
32
33 // handle by the kernel system call function
34 static int
35 sys_exit(uint64_t arg[]) {
36     int error_code = (int)arg[0];
37     return do_exit(error_code);
38 }
39

```

4. `do_exit()` recycles the user process

- Set the state of current process to be `PROC_ZOMBIE`. Waiting for the parent process to recycle the process. Set the `current->exit_code = error_code`; means that this process can not be scheduled
- If the parent of the current process is waiting for a child, wake up the parent to reclaim resources

- Then start the interrupt, execute the `schedule()` function, and select a new process to execute

```
1 // do_exit - called by sys_exit
2 // 1. call exit_mmap & put_pgdir & mm_destroy to free the almost all memory space of
   process
3 // 2. set process' state as PROC_ZOMBIE, then call wakeup_proc(parent) to ask parent
   reclaim itself.
4 // 3. call scheduler to switch to other process
5 int
6 do_exit(int error_code) {
7     if (current == idleproc) {
8         panic("idleproc exit.\n");
9     }
10    if (current == initproc) {
11        panic("initproc exit.\n");
12    }
13    struct mm_struct *mm = current->mm;
14    if (mm != NULL) {
15        lcr3(boot_cr3);
16        if (mm_count_dec(mm) == 0) {
17            exit_mmap(mm);
18            put_pgdir(mm);
19            mm_destroy(mm);
20        }
21        current->mm = NULL;
22    }
23    current->state = PROC_ZOMBIE;
24    current->exit_code = error_code;
25    bool intr_flag;
26    struct proc_struct *proc;
27    local_intr_save(intr_flag);
28    {
29        proc = current->parent;
30        if (proc->wait_state == WT_CHILD) {
31            wakeup_proc(proc);
32        }
33        while (current->cptr != NULL) {
34            proc = current->cptr;
35            current->cptr = proc->optr;
36
37            proc->yptr = NULL;
38            if ((proc->optr = initproc->cptr) != NULL) {
39                initproc->cptr->yptr = proc;
40            }
41            proc->parent = initproc;
42            initproc->cptr = proc;
43            if (proc->state == PROC_ZOMBIE) {
44                if (initproc->wait_state == WT_CHILD) {
45                    wakeup_proc(initproc);
46                }
47            }
48        }
49    }
50    local_intr_restore(intr_flag);
51    schedule();
52    panic("do_exit will not return!! %d.\n", current->pid);
53 }
```

- 1 首先执行 `lcr3(boot_cr3)`；切换到内核的页表上，这样用户进程就只能在内核的虚拟地址空间上执行，因为内核权限高。如果当前进程的被调用数减一后等于0，那么就没有其他进程在使用了，就可以进行回收，先回收内存资源，调用 `exit_mmap` 函数释放 `mm` 中的 `vma` 描述的进程合法空间中实际分配的内存，然后把对应的页表项内容清空，最后把页表项和页目录表清空。然后调用 `put_pgdir` 函数释放页目录表所占用的内存。最后调用 `mm_destroy` 释放 `vma` 与 `mm` 的内存。把 `mm` 置为 `NULL`，表示与当前进程相关的用户虚拟内存空间和对应的内存管理成员变量所占的内核虚拟内存空间已经回收完毕

Q4

进程如何变成僵尸进程?

After child process finishes its execution. The state of child process will be set `PROC_ZOMBIE`. Indicates that the child process is in the exit state and requires current process (parent process) complete the final reclamation of the child process. If the current process is still in execution and has not been executed to `wait` to reclaim the child process, then this child process will become the zombie process.