

Lab 5: RV64 用户态程序

课程名称：操作系统 指导老师：寿黎但

姓名：

庄毅非 3200105872 胡競文 3200105990

分工：

庄毅非：修改 head.S, start_kernel, 完成 elf 文件加载

胡競文：实现进程初始化, 实现 ecall 中断处理和系统调用

1、实验目的

- 创建用户态进程, 并设置 `sstatus` 来完成内核态转换至用户态。
- 正确设置用户进程的用户态栈和内核态栈, 并在异常处理时正确切换。
- 补充异常处理逻辑, 完成指定的系统调用 (`SYS_WRITE`, `SYS_GETPID`) 功能。

2、实验环境

- Environment in previous labs

3、实验步骤

3.1 准备工作

- 这一步主要就是将仓库中的 `elf.h` 等文件添加到 lab5, 并修改对应的 `makefile` 来确保项目可以编译。
- 修改根目录下的 `Makefile`, 在 `all` 和 `clean` 操作中添加对 `user` 文件夹中的 `Makefile` 文件的操作

```
1 all:
2     ${MAKE} -C lib all
3     ${MAKE} -C init all
4     ${MAKE} -C user all
5     ${MAKE} -C arch/riscv all
6     @echo -e '\n'Build Finished OK
7
8 run: all
9     @echo Launch the qemu .....
10    @qemu-system-riscv64 -nographic -machine virt -kernel vmlinux -bios
    default
11
12 debug: all
13    @echo Launch the qemu for debug .....
14    @qemu-system-riscv64 -nographic -machine virt -kernel vmlinux -bios
    default -S -s
15
16 clean:
17    ${MAKE} -C lib clean
```

```

18     ${MAKE} -C init clean
19     ${MAKE} -C user clean
20     ${MAKE} -C arch/riscv clean
21     $(shell test -f vmlinux && rm vmlinux)
22     $(shell test -f System.map && rm System.map)
23     @echo -e '\n'Clean Finished
24

```

3.2 创建用户态进程

- 修改 `proc.h` 文件夹中的 `NR_TASK` 为 `1 + 4`，表示实验中只有一个 `idle` 内核线程和 4 个用户线程

```

1  #define NR_TASKS (1 + 4)

```

- 修改 `thread_struct` 如下，添加对 `sepc`, `sstatus` 和 `sscratch` 寄存器的存储，在 `__switch_ro` 中使用

```

1  struct thread_struct
2  {
3      unsigned long ra; // return address
4      unsigned long sp;
5      unsigned long s[12];
6      unsigned long sepc;
7      unsigned long sstatus;
8      unsigned long sscratch;
9  };
10

```

- 修改 `task_init` 如下

```

1  // kernel/vm.c
2  // 将 src 开始的一个 page 拷贝到 dest 为起始的一个 page 上
3  void copy(void *src, void *dest)
4  {
5      unsigned long *beh = (unsigned long *) (dest);
6      unsigned long *pre = (unsigned long *) (src);
7      for (int i = 0; i < PGSIZE / 8; i++)
8      {
9          ((unsigned long *) (beh))[i] = ((unsigned long *) (pre))[i];
10     }
11 }
12 // kernel/proc.c
13 void task_init()
14 {

```

```

15 // 初始化 idle 线程
16 idle = (struct task_struct *)kalloc();
17 idle->state = TASK_RUNNING;
18 idle->counter = 0;
19 idle->priority = 0;
20 idle->pid = 0;
21 current = idle;
22 task[0] = idle;
23
24 for (int i = 1; i ≤ NR_TASKS - 1; ++i)
25 {
26     task[i] = (struct task_struct *)kalloc();
27     task[i]->state = TASK_RUNNING;
28     task[i]->counter = 0;
29     task[i]->priority = rand();
30     task[i]->pid = i;
31     task[i]->thread.ra = (unsigned long)__dummy;
32     task[i]->thread.sp = (unsigned long)task[i] + PGSIZE;
33     task[i]->thread_info = (struct thread_info *) (alloc_page());
34     task[i]->thread_info->kernel_sp = task[i]->thread.sp;
35     // 分配用户栈
36     task[i]->thread_info->user_sp = (alloc_page());
37     // 用户页表初始化
38     pagetable_t pgtbl = (unsigned long *) (alloc_page());
39     // 将内核页表拷贝到用户页表中
40     copy(swapper_pg_dir, pgtbl);
41     // 完成用户栈的映射
42     create_mapping(pgtbl, USER_END - PGSIZE, (unsigned long)(task[i]-
>thread_info->user_sp - PA2VA_OFFSET), PGSIZE, READABLE | WRITABLE | USER |
VALID);
43     // 用户程序占据的页面数
44     unsigned long user_page_count = PGROUNDUP(uapp_end - uapp_start) /
4096;
45     // 为用户程序分配所需的 pages
46     unsigned long *user_app_pages = (unsigned long
*)alloc_pages(user_page_count);
47     for (int i = 0; i < user_page_count; i++)
48     {
49         // 将用户程序拷贝到各自维护的 pages 中, 保障多个进程之间的地址空间隔离
50         copy((unsigned long *)((unsigned long)(uapp_start) + PGSIZE * i),
(unsigned long *)((unsigned long)(user_app_pages) + PGSIZE * i));
51     }
52     // 完成用户代码段的映射

```

```

53     create_mapping(pgtbl, USER_START, (unsigned long)((unsigned long)
(user_app_pages)-PA2VA_OFFSET), (unsigned long)uapp_end - (unsigned
long)uapp_start, USER | VALID | READABLE | WRITABLE | EXECUTABLE);
54     // set sepc
55     task[i]→thread.sepc = (unsigned long)(USER_START);
56     // set sstatus
57     unsigned long sstatus = csr_read(sstatus);
58     // set spp bit to 0
59     sstatus &= ~(1 << 8);
60     // set SPIE bit to 1
61     sstatus |= (1 << 5);
62     // set SUM bit to 1
63     sstatus |= (1 << 18);
64     task[i]→thread.sstatus = sstatus;
65     // user mode stack
66     task[i]→thread.sscratch = USER_END;
67     // set thread satp
68     unsigned long satp = csr_read(satp);
69     task[i]→pgd = (pagetable_t)(((satp >> 44) << 44) | (((unsigned long)
(pgtbl)-PA2VA_OFFSET) >> 12));
70 }
71 /* YOUR CODE HERE */
72 printk("...proc_init done!\n");
73 }
74

```

- 修改 `__switch_to`, 加入保存 / 恢复 `sepc` `sstatus` `sscratch` 以及 切换页表的逻辑。

```

1  .globl __switch_to
2  __switch_to:
3      # save state to prev process
4      # YOUR CODE HERE
5      # skip the first 5 virables
6      addi t0, a0, 8 * 5
7      sd ra, 0(t0)
8      sd sp, 8(t0)
9
10     addi t0, t0, 8 * 2
11     sd s0, 0*8(t0)
12     sd s1, 1*8(t0)
13     sd s2, 2*8(t0)
14     sd s3, 3*8(t0)
15     sd s4, 4*8(t0)
16     sd s5, 5*8(t0)
17     sd s6, 6*8(t0)

```

```

18     sd s7, 7*8(t0)
19     sd s8, 8*8(t0)
20     sd s9, 9*8(t0)
21     sd s10, 10*8(t0)
22     sd s11, 11*8(t0)
23     # save sepc, sstatus and sscratch
24     addi t0,t0,12 * 8
25     csrr s0,sepc
26     sd s0, 0 * 8(t0)
27     csrr s0,sstatus
28     sd s0,1 * 8(t0)
29     csrr s0,sscratch
30     sd s0, 2 * 8(t0)
31
32     addi t0,t0,3 * 8
33     # save satp
34     csrr s0,satp
35     sd s0,0(t0)
36     # restore state from next process
37     # YOUR CODE HERE
38     addi t1, a1, 8 * 5
39     ld ra, 0(t1)
40     ld sp, 8(t1)
41     addi t1, t1, 8 * 2
42     ld s0, 0*8(t1)
43     ld s1, 1*8(t1)
44     ld s2, 2*8(t1)
45     ld s3, 3*8(t1)
46     ld s4, 4*8(t1)
47     ld s5, 5*8(t1)
48     ld s6, 6*8(t1)
49     ld s7, 7*8(t1)
50     ld s8, 8*8(t1)
51     ld s9, 9*8(t1)
52     ld s10, 10*8(t1)
53     ld s11, 11*8(t1)
54     # restore sepc, sstatus and sscratch
55     addi t1,t1,12 * 8
56     ld t2, 0(t1)
57     csrw sepc, t2
58     ld t2, 8(t1)
59     csrw sstatus, t2
60     ld t2, 16(t1)
61     csrw sscratch, t2
62     addi t1,t1,3 * 8;

```

```

63     ld t2,0(t1)
64     # switch satp
65     csrw satp,t2
66     # flush tlb and lcache
67     sfence.vma
68     fence.i
69     ret
70

```

3.3 修改中断入口 / 返回逻辑 (`_trap`) 以及中断处理函数 (`trap_handler`)

- 修改 `__dummy`

```

1  __dummy:
2      # 在返回用户段的时候, 交换 sp 寄存器和 sscratch 寄存器, 完成用户栈和内核栈之间的切换
3      csrrw sp, sscratch, sp
4      sret

```

- 修改 `_traps` , 如果触发中断的是用户线程 (根据 `sscratch` 寄存器是否为 0 判断) , 那么执行栈的切换

```

1  #在本实验中, 触发中断的线程应该只有用户线程, 所以这里只留下 csrrw sp,sscratch,sp 指令,
    #程序也可以正常运行
2      csrrw sp,sscratch,sp
3      addi sp,sp,-34 * 8
4
5      sd x1, 1 * 8(sp)
6      addi x1,sp,34 * 8
7      sd x1, 2 * 8(sp)
8      sd x3, 3 * 8(sp)
9      ...
10     ld sp,2 * 8(sp)
11     #交换内核栈和用户栈
12     csrrw sp,sscratch,sp
13     sret

```

- 修改 `trap_handler` , 完成对 `ecall` 的捕获 (如果 `scause` 的第八位为 1, 首位为 0, 表示用户调用 `ecall`)

```

1  // syscall.h
2  struct pt_regs
3  {
4      // unsigned long x[32];
5      unsigned long zero;

```

```

6   unsigned long ra;
7   unsigned long sp;
8   unsigned long gp;
9   unsigned long tp;
10  unsigned long t0, t1, t2;
11  unsigned long fp;
12  unsigned long s1;
13  unsigned long a0, a1, a2, a3, a4, a5, a6, a7;
14  unsigned long s2, s3, s4, s5, s6, s7, s8, s9, s10, s11;
15  unsigned long t3, t4, t5, t6;
16  unsigned long sepc;
17  unsigned long sstatus;
18 };
19 // trap.c
20 void trap_handler(unsigned long scause, unsigned long sepc, struct pt_regs
*regs)
21 {
22     if ((scause >> 63) & 1)
23     {
24         // interrupt
25         if ((scause & 5) == 5)
26         {
27             clock_set_next_event();
28             do_timer();
29         }
30     }
31     else
32     {
33         // exception
34         if (scause & 8)
35         {
36             // ecall from u_mode
37             syscall(regs);
38         }
39     }
40 }

```

3.4 添加系统调用

- 添加 `syscall.h` 文件

```

1  #ifndef __SYSCALL_H__
2  #define __SYSCALL_H__
3

```

```

4  #define SYS_WRITE 64
5  #define SYS_GETPID 172
6
7  struct pt_regs
8  {
9      // unsigned long x[32];
10     unsigned long zero;
11     unsigned long ra;
12     unsigned long sp;
13     unsigned long gp;
14     unsigned long tp;
15     unsigned long t0, t1, t2;
16     unsigned long fp;
17     unsigned long s1;
18     unsigned long a0, a1, a2, a3, a4, a5, a6, a7;
19     unsigned long s2, s3, s4, s5, s6, s7, s8, s9, s10, s11;
20     unsigned long t3, t4, t5, t6;
21     unsigned long sepc;
22     unsigned long sstatus;
23 };
24 int sys_write(unsigned int fd, const char *buf, unsigned long count);
25 int sys_getpid();
26 void syscall(struct pt_regs *);
27 #endif

```

- `syscall.c` 文件内容如下

```

1  #include "syscall.h"
2  #include "printk.h"
3  #include "mm.h"
4  #include "defs.h"
5  #include "vm.h"
6  #include "proc.h"
7  extern struct task_struct *current;
8  int sys_write(unsigned int fd, const char *buf, unsigned long count)
9  {
10     // 分配新的连续 page 存储用户要输出的字符串
11     char *buffer = (char *) (alloc_page(PGROUNDDUP(count) / PGSIZE));
12     // 拷贝用户字符串
13     for (int i = 0; i < count; i++)
14     {
15         buffer[i] = buf[i];
16     }
17     // 设置结束字符
18     buffer[count] = '\0';

```



```

19     // 调用 printk 函数完成输出, 记录输出的字符数
20     int result = printk(buffer);
21     // 释放申请的内存页面
22     free_pages((unsigned long)(buffer));
23     return result;
24 }
25 int sys_getpid()
26 {
27     // 返回当前线程的线程号
28     return current->pid;
29 }
30 void syscall(struct pt_regs *regs)
31 {
32     unsigned long sys_code = regs->a7;
33     int return_value = -1;
34     // 判断服务类型
35     switch (sys_code)
36     {
37     case SYS_GETPID:
38     {
39         return_value = sys_getpid();
40         break;
41     }
42     case SYS_WRITE:
43     {
44         unsigned long fd = regs->a0;
45         unsigned long buf = regs->a1;
46         unsigned long count = regs->a2;
47         return_value = sys_write(fd, (const char *)buf, count);
48         break;
49     }
50     }
51     // write result to a0
52     regs->a0 = return_value;
53     regs->sepc += 4;
54 }

```

3.5 修改 head.S 以及 start_kernel

- 修改 `head.S`

```

1  # 将原来对 sstatus 设置的 sie bit 代码注释, 保证线程调度的时候不会被中断
2  # set sstatus[SIE] = 1
3  # csrr t0, sstatus
4  # ori t0, t0, 0x2
5  # csw sstatus, t0

```

- 修改 start_kernel 函数

```

1  #include "printk.h"
2  #include "sbi.h"
3  #include "proc.h"
4  extern void test();
5
6  int start_kernel()
7  {
8      // 调用 schedule 函数, 直接触发调度
9      schedule();
10     test();
11
12     return 0;
13 }
14

```

3.6 测试二进制文件能否运行

- 程序运行截图如下

```

[U-MODE] pid: 4, sp is 0000003fffffffffe0, this is print No.1
[U-MODE] pid: 3, sp is 0000003fffffffffe0, this is print No.1
[U-MODE] pid: 1, sp is 0000003fffffffffe0, this is print No.1
[U-MODE] pid: 1, sp is 0000003fffffffffe0, this is print No.2
[U-MODE] pid: 2, sp is 0000003fffffffffe0, this is print No.1
[U-MODE] pid: 2, sp is 0000003fffffffffe0, this is print No.2
[U-MODE] pid: 3, sp is 0000003fffffffffe0, this is print No.2
[U-MODE] pid: 1, sp is 0000003fffffffffe0, this is print No.3
[U-MODE] pid: 4, sp is 0000003fffffffffe0, this is print No.2
[U-MODE] pid: 4, sp is 0000003fffffffffe0, this is print No.3
[U-MODE] pid: 2, sp is 0000003fffffffffe0, this is print No.3

```

3.7 添加 ELF 支持

- 将 `user/uapp.S` 修改如下，使程序在运行的时候载入文件从 `uapp.bin` 变为 `uapp.elf`

```
1  /* user/uapp.S */
2  .section .uapp
3
4  .incbin "uapp"
```

- 修改线程初始化代码如下

```
1  static void load_program(struct task_struct *task)
2  {
3      Elf64_Ehdr *ehdr = (Elf64_Ehdr *)uapp_start;
4
5      unsigned long phdr_start = (unsigned long)ehdr + ehdr->e_phoff;
6      int phdr_cnt = ehdr->e_phnum;
7      pagetable_t pgtbl = (unsigned long *)((((unsigned long)(task->pgd)) &
0xffffffff) << 12) + PA2VA_OFFSET);
8      Elf64_Phdr *phdr;
9      int load_phdr_cnt = 0;
10     for (int i = 0; i < phdr_cnt; i++)
11     {
12         phdr = (Elf64_Phdr *) (phdr_start + sizeof(Elf64_Phdr) * i);
13         if (phdr->p_type == PT_LOAD)
14         {
15             // do mapping
16             // code...
17             // 读取载入段虚拟地址
18             unsigned long va = phdr->p_vaddr;
19             // 计算映射段在文件对应的大小需要多少个 page
20             unsigned long filePages = PGROUNDUP(phdr->p_filesz) / PGSIZE;
21             // 计算映射段在内存中对应的大小需要多少个 page
22             unsigned long memPages = (PGROUNDUP(phdr->p_memsz)) / PGSIZE;
23             // 复制代码段到本进程维护的页表中，保证不同进程之间使用的地址空间是隔离的
24             unsigned long mappedPage = alloc_pages(memPages);
25             for (int i = 0; i < PGROUNDUP(phdr->p_filesz) / 8; i++)
26             {
27                 ((unsigned long *) (mappedPage))[i] = ((unsigned long *)
(PGROUNDDOWN(((unsigned long)uapp_start + phdr->p_offset))))[i];
28             }
29             // 将上述映射的文件长度对应的 page 映射到页表中
30             create_mapping(pgtbl, va, (unsigned long)(mappedPage)-
PA2VA_OFFSET, phdr->p_filesz, (phdr->p_flags << 1) | USER | VALID);
31             // 将多余的内存长度对应的 page 映射到页表中
```

```

32         create_mapping(pgtbl, va + PGSIZE * filePages, (unsigned long)
(mappedPage) + PGSIZE * filePages - PA2VA_OFFSET, (memPages - filePages) *
PGSIZE, (phdr->p_flags << 1) | USER | VALID);
33         load_phdr_cnt++;
34     }
35 }
36 // 设置 sepc 为 elf 文件的起始代码
37 task->thread.sepc = ehdr->e_entry;
38 }
39 void task_init()
40 {
41     idle = (struct task_struct *)kalloc();
42     idle->state = TASK_RUNNING;
43     idle->counter = 0;
44     idle->priority = 0;
45     idle->pid = 0;
46     current = idle;
47     task[0] = idle;
48
49     for (int i = 1; i ≤ NR_TASKS - 1; ++i)
50     {
51         task[i] = (struct task_struct *)kalloc();
52         task[i]->state = TASK_RUNNING;
53         task[i]->counter = 0;
54         task[i]->priority = rand();
55         task[i]->pid = i;
56         task[i]->thread.ra = (unsigned long)__dummy;
57         task[i]->thread.sp = (unsigned long)task[i] + PGSIZE;
58         task[i]->thread_info = (struct thread_info *) (alloc_page());
59         task[i]->thread_info->kernel_sp = task[i]->thread.sp;
60         task[i]->thread_info->user_sp = (alloc_page());
61         // 用户页表初始化
62         pagetable_t pgtbl = (unsigned long *) (alloc_page());
63         // 拷贝内核页表
64         copy(swapper_pg_dir, pgtbl);
65         // 映射用户栈
66         create_mapping(pgtbl, USER_END - PGSIZE, (unsigned long)(task[i]-
>thread_info->user_sp - PA2VA_OFFSET), PGSIZE, READABLE | WRITABLE | USER |
VALID);
67         // set sstatus
68         unsigned long sstatus = csr_read(sstatus);
69         // set spp bit to 0
70         sstatus &= (~(1 << 8));
71         // set SPIE bit to 1
72         sstatus |= (1 << 5);

```

```

73         // set SUM bit to 1
74         sstatus |= (1 << 18);
75         task[i]→thread.sstatus = sstatus;
76         // user mode stack
77         task[i]→thread.sscratch = USER_END;
78
79         // printk("SET [PID = %d] COUNTER = %d\n", task[i] → pid, task[i] →
counter);
80         unsigned long satp = csr_read(satp);
81         task[i]→pgd = (pagetable_t)((((satp >> 44) << 44) | (((unsigned long)
(pgtbl)-PA2VA_OFFSET) >> 12)));
82         // load uapp.elf
83         load_program(task[i]);
84     }
85     /* YOUR CODE HERE */
86     printk("...proc_init done!\n");
87 }
88

```

- 完成上述修改之后，每一个用户线程在初始化的时候都会将 elf 文件对应的 load 段装载到自己维护的页表中，这样实现了不同进程之间的地址空间的隔离，程序也能够运行 elf 文件中对应的代码。
- 程序运行输出截图如下。

```

Boot HART ID : 0
Boot HART Domain : root
Boot HART ISA : rv64imafdcsh
Boot HART Features : scounteren,mcounteren,mcountinhibit,time
Boot HART PMP Count : 16
Boot HART PMP Granularity : 4
Boot HART PMP Address Bits: 54
Boot HART MHPM Count : 16
Boot HART MIDELEG : 0x00000000000001666
Boot HART MEDELEG : 0x00000000000f0b509
...buddy_init done!
setup_vm_final finish
...proc_init done!
[U-MODE] pid: 4, sp is 0000003ffffffe0, this is print No.1
[U-MODE] pid: 3, sp is 0000003ffffffe0, this is print No.1
[U-MODE] pid: 1, sp is 0000003ffffffe0, this is print No.1
[U-MODE] pid: 1, sp is 0000003ffffffe0, this is print No.2
[U-MODE] pid: 2, sp is 0000003ffffffe0, this is print No.1
[U-MODE] pid: 2, sp is 0000003ffffffe0, this is print No.2
[U-MODE] pid: 3, sp is 0000003ffffffe0, this is print No.2
[U-MODE] pid: 1, sp is 0000003ffffffe0, this is print No.3
[U-MODE] pid: 4, sp is 0000003ffffffe0, this is print No.2
[U-MODE] pid: 1, sp is 0000003ffffffe0, this is print No.4
[U-MODE] pid: 4, sp is 0000003ffffffe0, this is print No.3
[U-MODE] pid: 2, sp is 0000003ffffffe0, this is print No.3
[U-MODE] pid: 2, sp is 0000003ffffffe0, this is print No.4
[U-MODE] pid: 2, sp is 0000003ffffffe0, this is print No.5
[U-MODE] pid: 3, sp is 0000003ffffffe0, this is print No.3

```

4. 思考题

4.1 我们在实验中使用的用户态线程和内核态线程的对应关系是怎样的？（一对一，一对多，多对一还是多对多）

程序中使用的用户态线程和内核态线程是 1 对 1 的关系。

4.2 为什么 Phdr 中，`p_filesz` 和 `p_memsz` 是不一样大的？

因为在 elf 文件中的一个 load 段可能有无需初始化的变量或者数组等数据，这些数据都会在程序装载的时候初始化为 0，如果在 elf 文件中为这些数据预留空间的话，就会造成磁盘空间的浪费。但是这些数据在内存中是需要对应的空间的，所以在许多情况下 `p_memsz` 会稍大于 `p_filesz`。

在 linuxbase.org 的文档中也有对 load 段 `p_memsz` 稍大于 `p_filesz` 段的说明。

PT_LOAD

The array element specifies a loadable segment, described by `p_filesz` and `p_memsz`. The bytes from the file are mapped to the beginning of the memory segment. If the segment's memory size (`p_memsz`) is larger than the file size (`p_filesz`), the "extra" bytes are defined to hold the value 0 and to follow the segment's initialized area. The file size may not be larger than the memory size. Loadable segment entries in the program header table appear in ascending order, sorted on the `p_vaddr` member.