

# Lab 4: RV64 用户态程序

3210105944 黄锦骏

## 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

- Same as previous labs.

## 3 实验步骤

### 3.1 准备工程

- 此次实验基于 lab4 同学所实现的代码进行。
- 需要修改 `vmlinux.ld.s`，将用户态程序 `uapp` 加载至 `.data` 段。按如下修改：

```
1  ...
2
3  .data : ALIGN(0x1000){
4      _sdata = .;
5
6      *(.sdata .sdata*)
7      *(.data .data.*)
8
9      _edata = .;
10
11     . = ALIGN(0x1000);
12     uapp_start = .;
13     *(.uapp .uapp*)
14     uapp_end = .;
15     . = ALIGN(0x1000);
16
17     } >ramv AT>ram
18
19  ...
```

如果你要使用 `uapp_start` 这个符号，可以在代码里这样来声明它： `extern char uapp_start[]`，这样就可以像一个字符数组一样来访问这块内存的内容。例如，程序的第一个字节就是 `uapp_start[0]`。

- 需要修改 `defs.h`，在 `defs.h` 添加 如下内容：

```
1  #define USER_START (0x0000000000000000) // user space start virtual address
2  #define USER_END   (0x0000004000000000) // user space end virtual address
```

- 从 `repo` 同步以下文件和文件夹。并按照下面的位置来放置这些新文件。值得注意的是，我们在 `mm` 中添加了 `buddy system`，但是也保证了原来调用的 `kalloc` 和 `kfree` 的兼容。你应该无需修改原先使用了 `kalloc` 的相关代码，如果出现兼容性问题可以联系助教。为了减小大家的工作量，我们替大家实现了 `Buddy System`，大家可以直接使用这些函数来管理内存：

```
1 // 分配 page_cnt 个页的地址空间，返回分配内存的地址。保证分配的内存存在虚拟地址和物理地址上
  都是连续的
2 uint64_t alloc_pages(uint64_t page_cnt);
3 // 相当于 alloc_pages(1);
4 uint64_t alloc_page();
5 // 释放从 addr 开始的之前按分配的内存
6 void free_pages(uint64_t addr);
```

```
1 .
2 |— arch
3 |   |— riscv
4 |       |— Makefile
5 |       |— include
6 |           |— mm.h
7 |           |— stdint.h
8 |           |— kernel
9 |           |— mm.c
10 |— include
11 |   |— elf.h (this is copied from newlib)
12 |— user
13     |— Makefile
14     |— getpid.c
15     |— link.lds
16     |— printf.c
17     |— start.s
18     |— stddef.h
19     |— stdio.h
20     |— syscall.h
21     |— uapp.s
```

- 修改根目录下的 `Makefile`, 将 `user` 纳入工程管理。

```
1 .PHONY:all run debug clean
2 all: clean
3     ...
4     ${MAKE} -C user all
5     ...
6
7 TEST:
8     ...
9     ${MAKE} -C user all
10    ...
11    ...
12
13 clean:
14    ...
15    ${MAKE} -C user clean
16    ...
```

- 在根目录下 `make` 会生成 `user/uapp.o` `user/uapp.elf` `user/uapp.bin`，以及我们最终测试使用的 ELF 可执行文件 `user/uapp`。通过 `objdump` 我们可以看到 `uapp` 使用 `ecall` 来调用 `SYSCALL` (在 U-Mode 下使用 `ecall` 会触发 `environment-call-from-U-mode` 异常)。从而将控制权交给处在 S-Mode 的 OS，由内核来处理相关异常。
- 在本次实验中，我们首先会将用户态程序 `strip` 成纯二进制文件来运行。这种情况下，用户程序运行的第一条指令位于二进制文件的开始位置，也就是说 `uapp_start` 处的指令就是我们要执行的第一条指令。我们将运行纯二进制文件作为第一步，在确认用户态的纯二进制文件能够运行后，我们再将存储到内存中的用户程序文件换为 ELF 来进行执行。

```

1  0000000000000004 <getpid>:
2      4:   fe010113          addi    sp,sp,-32
3      8:   00813c23          sd      s0,24(sp)
4      c:   02010413          addi    s0,sp,32
5     10:   fe843783          ld      a5,-24(s0)
6     14:   0ac00893          li      a7,172
7     18:   00000073          ecall                                <-
      SYS_GETPID
8     ...
9
10 00000000000000d8 <vprintfmt>:
11 ...
12 60c:   00070513          mv      a0,a4
13 610:   00068593          mv      a1,a3
14 614:   00060613          mv      a2,a2
15 618:   00000073          ecall                                <- SYS_WRITE
16 ...

```

## 3.2 创建用户态进程

### 1. 修改 `proc.h`

- 将 `NR_TASKS` 修改为 5
- 由于创建用户态进程要对 `sepc` `sstatus` `sscratch` 做设置，我们将其加入 `thread_struct` 中
- 由于多个用户态进程需要保证相对隔离，因此不可以共用页表。我们为每个用户态进程都创建一个页表 `pgd`。

```

1  #ifndef _PROC_H
2  #define _PROC_H
3  #include "types.h"
4
5  #define NR_TASKS (1+4) // 用于控制 最大线程数量 (idle 线程 + 31 内核线程)

```

```

6
7 typedef unsigned long* pagetable_t;
8
9 /* 线程状态段数据结构 */
10 struct thread_struct {
11     uint64 ra;
12     uint64 sp;
13     uint64 s[12];
14
15     uint64 sepc, sstatus, sscratch;
16 };
17
18 /* 线程数据结构 */
19 struct task_struct {
20     struct thread_info thread_info;
21     uint64 state; // 线程状态
22     uint64 counter; // 运行剩余时间
23     uint64 priority; // 运行优先级 1最低 10最高
24     uint64 pid; // 线程id
25
26     struct thread_struct thread;
27
28     pagetable_t pgd;
29 };
30

```

## 2. 修改task\_init

- 对每个用户态进程，其拥有两个 stack：U-Mode Stack 以及 S-Mode Stack，其中 S-Mode stack 在 lab3 中我们已经设置好了。我们可以通过 alloc\_page 接口申请一个空的页面来作为 U-Mode Stack。
- 为每个用户态进程创建自己的页表 并将 uapp 所在页面，以及 U-Mode Stack 做相应的映射，同时为了避免 U-Mode 和 S-Mode 切换的时候切换页表，我们也将内核页表（swapper\_pg\_dir）复制到每个进程的页表中。注意程序运行过程中，有部分数据不在栈上，而在初始化的过程中就已经被分配了空间（比如我们的 uapp 中的 counter 变量），所以二进制文件需要先被拷贝到一块某个进程专用的内存之后再行映射，防止所有的进程共享数据，造成期望外的进程间相互影响。

```

1 // /arch/riscv/kernel/proc.c
2 //分配页表内存
3 task[i]->pgd = (pagetable_t)kalloc();
4 memset(task[i]->pgd, 0, PGSIZE);
5 // 拷贝内核页表到进程页表
6 for(int j = 0; j < PGSIZE / 8; j++){
7     task[i]->pgd[j]=swapper_pg_dir[j];
8 }
9 // 申请uapp在进程内的专用内存
10 uint64 uapp = alloc_pages((PGROUNDUP((uint64)uapp_end) -
    (uint64)uapp_start) / PGSIZE);
11 //将uapp.bin拷贝到进程内存
12 for(int t = 0; t < (PGROUNDUP((uint64)uapp_end) -
    (uint64)uapp_start) / PGSIZE; t++){

```

```

13     uint64* src = (uint64*)(uapp_start + t * PGSIZE);
14     uint64* dst = (uint64*)(uapp + t * PGSIZE);
15     for (int k = 0; k < PGSIZE / 8; k++) {
16         dst[k] = src[k];
17     }
18 }
19 //创建uapp的映射
20 create_mapping(task[i]->pgd, (uint64)USER_START, (uint64)uapp -
    (uint64)PA2VA_OFFSET, (uint64)uapp_end - (uint64)uapp_start, 31);
21 //创建用户栈的映射, 映射到用户空间的最后一个页面
22 create_mapping(task[i]->pgd, (uint64)USER_END - (uint64)PGSIZE,
    (uint64)kalloc() - (uint64)PA2VA_OFFSET, (uint64)PGSIZE, 23);

```

- 对每个用户态进程我们需要将 `sepc` 修改为 `USER_START`, 配置修改好 `sstatus` 中的 `SPP` (使得 `sret` 返回至 U-Mode), `SPIE` (`sret` 之后开启中断), `SUM` (S-Mode 可以访问 User 页面), `sscratch` 设置为 U-Mode 的 `sp`, 其值为 `USER_END` (即 U-Mode `stack` 被放置在 `user space` 的最后一个页面)。

```

1 // /arch/riscv/kernel/proc.c
2 task[i]->thread.sepc = USER_START;
3 task[i]->thread.sstatus = (1ULL << 5) + (1ULL << 18);
4 task[i]->thread.sscratch = USER_END;

```

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

```

1 # /arch/riscv/kernel/entry.S
2 __switch_to:
3     ...
4
5     # save state of sepc, sstatus, sscratch
6     csrr t1, sepc
7     sd t1, 160(a0)
8     csrr t1, sstatus
9     sd t1, 168(a0)
10    csrr t1, sscratch
11    sd t1, 176(a0)
12
13    ...
14
15    # load state of sepc, sstatus, sscratch
16    ld t1, 160(a1)
17    csrwr sepc, t1
18    ld t1, 168(a1)
19    csrwr sstatus, t1
20    ld t1, 176(a1)
21    csrwr sscratch, t1
22
23    #切换页表
24    li t5, PA2VA_OFFSET
25    li t0, 8
26    slli t0, t0, 60 # Mode = 8
27    ld t1, 184(a1)
28    sub t1, t1, t5
29    srli t1, t1, 12

```

```

30     or t0, t0, t1
31
32     csrw satp, t0
33     sfence.vma zero, zero
34     fence.i
35     ret

```

### 3.3 修改中断入口/返回逻辑 ( \_trap ) 以及中断处理函数 ( trap\_handler )

- 与 ARM 架构不同的是, RISC-V 中只有一个栈指针寄存器( sp ), 因此需要我们来完成用户栈与内核栈的切换。
- 由于我们的用户态进程运行在 U-Mode 下, 使用的运行栈也是 U-Mode Stack, 因此当触发异常时, 我们首先要对栈进行切换 ( U-Mode Stack -> S-Mode Stack )。同理 让我们完成了异常处理, 从 S-Mode 返回至 U-Mode, 也需要进行栈切换 ( S-Mode Stack -> U-Mode Stack )。

1. 修改 \_\_dummy。在 4.2 中我们初始化时, thread\_struct.sp 保存了 S-Mode sp, thread\_struct.sscratch 保存了 U-Mode sp, 因此在 S-Mode -> U->Mode 的时候, 我们只需要交换对应的寄存器的值即可。

```

1  __dummy:
2      # la t0, dummy # 将 dummy() 的地址加载到寄存器 t0 中
3      # csrw sepc, t0 # 将 dummy() 的地址写入寄存器 sepc
4
5      # 交换 u-mode stack和 s-mode stack
6      csrr t1, sscratch
7      csrw sscratch, sp
8      add sp, t1, x0
9
10     sret

```

2. 修改 \_trap。同理在 \_trap 的首尾我们都需要做类似的操作。注意如果是 内核线程( 没有 U-Mode Stack )触发了异常, 则不需要进行切换。(内核线程的 sp 永远指向的 S-Mode Stack, sscratch 为 0)

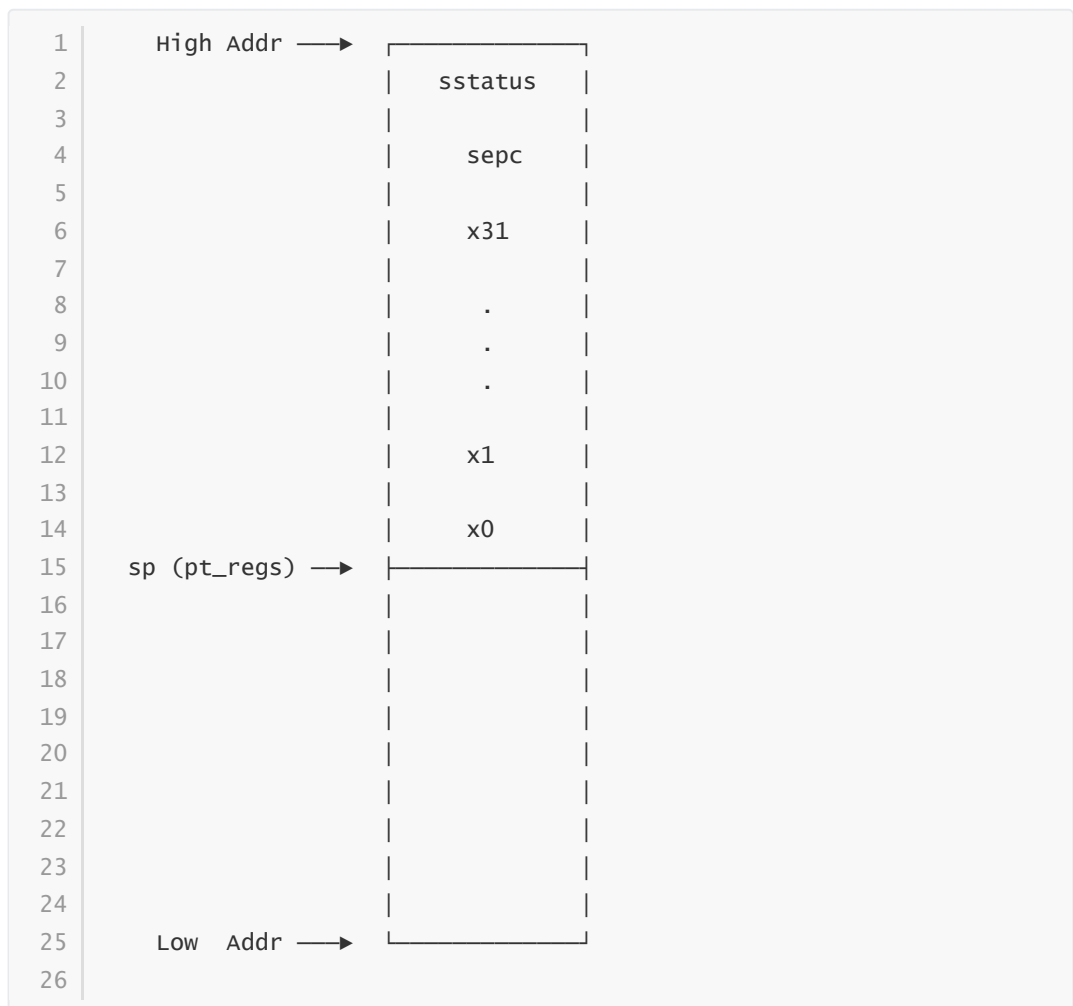
```

1  _traps:
2      # -----
3      #判断为内核线程 / 用户线程
4      csrr t1, sscratch
5      beq t1, x0, _trap_core
6
7      csrw sscratch, sp
8      add sp, t1, x0
9
10     ...
11
12     csrr t1, sscratch
13     beq t1, x0, _trap_end
14     csrw sscratch, sp
15     add sp, t1, x0
16
17     ...

```

3. 修改 `trap_handler`, `uapp` 使用 `ecall` 会产生 `ECALL_FROM_U_MODE` **exception**。因此我们需要在 `trap_handler` 里面进行捕获

- 经查询得知, `ECALL_FROM_U_MODE` exception对应的scause为 `0x8`, 捕获到就调用 `sys_call`
- 这里需要解释新增加的第三个参数 `regs`, 在 `_trap` 中我们将寄存器的内容连续的保存在 S-Mode Stack上, 因此我们可以将这一段看做一个叫做 `pt_regs` 的结构体。我们可以从这个结构体中取到相应的寄存器的值 (比如 `syscall` 中我们需要从 `a0 ~ a7` 寄存器中取到参数)。这个结构体中的值也可以按需添加, 同时需要在 `_trap` 中存入对应的寄存器值以供使用, 示例如下图:



- `trap_handler` 的实现以及 `struct pt_regs` 的定义如下:

```
1 struct pt_regs {
2     uint64 sstatus;
3     uint64 sepc;
4     uint64 x[32];
5 };
6
7
8 void trap_handler(uint64 scause, uint64 sepc, struct pt_regs *regs) {
9
10    if((scause << 63) && ((scause & 0x5) == 0x5))
11    {
12        // printk("[S] Supervisor Mode Timer Interrupt\n");
13        clock_set_next_event();
14        do_timer();
15    }
16 }
```

```

15     }
16     //Interrupt =1 && exception code == 8
17     else if((scause&0x8) == 8)
18     {
19         sys_call(regs);
20     }else{
21         printk("scause=%ld__",scause);
22     }
23     return;
24 }

```

### 3.4 添加系统调用

- 本次实验要求的系统调用函数原型以及具体功能如下：
  - 64 号系统调用 `sys_write(unsigned int fd, const char* buf, size_t count)` 该调用将用户态传递的字符串打印到屏幕上，此处fd为标准输出（1），buf为用户需要打印的起始地址，count为字符串长度，返回打印的字符数。（具体见 user/printf.c）
  - 172 号系统调用 `sys_getpid()` 该调用从current中获取当前的pid放入a0中返回，无参数。（具体见 user/getpid.c）
- 增加 `syscall.c` `syscall.h` 文件，并在其中实现 `getpid` 以及 `write` 逻辑。
- 系统调用的返回参数放置在 `a0` 中 (不可以直接修改寄存器，应该修改 `regs` 中保存的内容)。
- 实现的 `syscall` 如下：

```

1 void sys_call(struct pt_regs *regs){
2     uint64 func = regs->x[17];
3
4     if(func == 172){
5         regs->x[10] = current->pid;
6     }
7     else if(func == 64){
8         //fd(a0) = 1
9         if(regs->x[10] == 1){
10            //a1 = buf, a2 = count
11            ((char*)(regs->x[11]))[regs->x[12]] = '\0';
12            //a0 = 输出长度
13            printk((char*)(regs->x[11]));
14            regs->x[10] = regs->x[12];
15        }
16    }
17    return;
18 }

```

- 针对系统调用这一类异常，我们需要手动将 `sepc + 4`（`sepc` 记录的是触发异常的指令地址，由于系统调用这类异常处理完成之后，我们应该继续执行后续的指令，因此需要我们手动修改 `sepc` 的地址，使得 `sret` 之后程序继续执行）。



```
1 # /arch/riscv/kernel/entry.S
2     ...
3
4     ld t0, 8(x2)
5     addi t0, t0, 4 # 手动修改 sepc 的地址, 使得 sret 之后 程序继续执行
6     csrw sepc, t0
7
8     ...
```

### 3.5 修改 head.S 以及 start\_kernel

- 在之前的 lab 中，在 OS boot 之后，我们需要等待一个时间片，才会进行调度。我们现在更改为 OS boot 完成之后立即调度 uapp 运行。
- 在 start\_kernel 中调用 schedule() 注意放置在 test() 之前。

```
1 int start_kernel() {
2     printk("2022");
3     printk(" Hello RISC-V\n");
4
5     schedule();
6     test(); // DO NOT DELETE !!!
7
8     return 0;
9 }
10
```

- 将 head.S 中 enable interrupt sstatus.SIE 逻辑注释，确保 schedule 过程不受中断影响

```
1 # /arch/riscv/kernel/head.S
2
3 # 使schedule不受中断影响
4 #li t6, 0x2
5 #csrr t3, sstatus
6 #or t3, t3, t6
7 #csrw sstatus, t3
```

## 4.6 测试纯二进制文件

- 测试输出如下:

```

1  OpenSBI v0.9
2
3      ____          ____  ____
4  /  __ \          / ____|  _ \  __|
5  | |  | |  _ _   _ _ _ _ | (___| | |
6  | |  | | ' _ \ / _ \ ' _ \ ___ \ |  <  | |
7  | |  | | |_) | ___/ | | |___) | |_) | |  _
8  \___/| ._/ \___|_| | |___/|___/___|
9      | |
      |_|

```

```
10
11 Platform Name           : riscv-virtio,qemu
12 Platform Features       : timer,mfdeleg
13 Platform HART Count     : 1
14 Firmware Base           : 0x80000000
15 Firmware Size           : 100 KB
16 Runtime SBI Version     : 0.2
17
18 Domain0 Name            : root
19 Domain0 Boot HART       : 0
20 Domain0 HARTs           : 0*
21 Domain0 Region00        : 0x0000000080000000-0x000000008001ffff ()
22 Domain0 Region01        : 0x0000000000000000-0xffffffffffffffff
    (R,W,X)
23 Domain0 Next Address    : 0x0000000080200000
24 Domain0 Next Arg1       : 0x0000000087000000
25 Domain0 Next Mode       : S-mode
26 Domain0 SysReset       : yes
27
28 Boot HART ID            : 0
29 Boot HART Domain        : root
30 Boot HART ISA           : rv64imafdcsu
31 Boot HART Features      : scounteren,mcounteren,time
32 Boot HART PMP Count     : 16
33 Boot HART PMP Granularity : 4
34 Boot HART PMP Address Bits: 54
35 Boot HART MHPM Count    : 0
36 Boot HART MHPM Count    : 0
37 Boot HART MIDELEG       : 0x0000000000000222
38 Boot HART MEDELEG       : 0x000000000000b109
39 ...buddy_init done!
40 ...proc_init done!
41 2022 Hello RISC-V
42 switch to [PID = 1 COUNTER = 4]
43 [U-MODE] pid: 1, sp is 0000003fffffffffe0, this is print No.1
44 [U-MODE] pid: 1, sp is 0000003fffffffffe0, this is print No.2
45 [U-MODE] pid: 1, sp is 0000003fffffffffe0, this is print No.3
46 [U-MODE] pid: 1, sp is 0000003fffffffffe0, this is print No.4
47 switch to [PID = 4 COUNTER = 5]
48 [U-MODE] pid: 4, sp is 0000003fffffffffe0, this is print No.1
49 [U-MODE] pid: 4, sp is 0000003fffffffffe0, this is print No.2
50 [U-MODE] pid: 4, sp is 0000003fffffffffe0, this is print No.3
51 [U-MODE] pid: 4, sp is 0000003fffffffffe0, this is print No.4
52 [U-MODE] pid: 4, sp is 0000003fffffffffe0, this is print No.5
53 switch to [PID = 3 COUNTER = 8]
54 [U-MODE] pid: 3, sp is 0000003fffffffffe0, this is print No.1
55 [U-MODE] pid: 3, sp is 0000003fffffffffe0, this is print No.2
56 [U-MODE] pid: 3, sp is 0000003fffffffffe0, this is print No.3
57 [U-MODE] pid: 3, sp is 0000003fffffffffe0, this is print No.4
58 [U-MODE] pid: 3, sp is 0000003fffffffffe0, this is print No.5
59 [U-MODE] pid: 3, sp is 0000003fffffffffe0, this is print No.6
60 [U-MODE] pid: 3, sp is 0000003fffffffffe0, this is print No.7
61 [U-MODE] pid: 3, sp is 0000003fffffffffe0, this is print No.8
62 switch to [PID = 2 COUNTER = 9]
63 [U-MODE] pid: 2, sp is 0000003fffffffffe0, this is print No.1
```

```

64 [U-MODE] pid: 2, sp is 0000003ffffffffffe0, this is print No.2
65 [U-MODE] pid: 2, sp is 0000003ffffffffffe0, this is print No.3
66 [U-MODE] pid: 2, sp is 0000003ffffffffffe0, this is print No.4
67 [U-MODE] pid: 2, sp is 0000003ffffffffffe0, this is print No.5
68 [U-MODE] pid: 2, sp is 0000003ffffffffffe0, this is print No.6
69 [U-MODE] pid: 2, sp is 0000003ffffffffffe0, this is print No.7
70 [U-MODE] pid: 2, sp is 0000003ffffffffffe0, this is print No.8
71 [U-MODE] pid: 2, sp is 0000003ffffffffffe0, this is print No.9
72 switch to [PID = 1 COUNTER = 1]
73 [U-MODE] pid: 1, sp is 0000003ffffffffffe0, this is print No.5
74 switch to [PID = 2 COUNTER = 4]
75 [U-MODE] pid: 2, sp is 0000003ffffffffffe0, this is print No.10
76 [U-MODE] pid: 2, sp is 0000003ffffffffffe0, this is print No.11
77 [U-MODE] pid: 2, sp is 0000003ffffffffffe0, this is print No.12
78 [U-MODE] pid: 2, sp is 0000003ffffffffffe0, this is print No.13
79 switch to [PID = 4 COUNTER = 4]
80 [U-MODE] pid: 4, sp is 0000003ffffffffffe0, this is print No.6
81 [U-MODE] pid: 4, sp is 0000003ffffffffffe0, this is print No.7
82 [U-MODE] pid: 4, sp is 0000003ffffffffffe0, this is print No.8
83 [U-MODE] pid: 4, sp is 0000003ffffffffffe0, this is print No.9
84 [U-MODE] pid: 4, sp is 0000003ffffffffffe0, this is print No.10
85 [U-MODE] pid: 4, sp is 0000003ffffffffffe0, this is print No.11
86 switch to [PID = 3 COUNTER = 5]
87 [U-MODE] pid: 3, sp is 0000003ffffffffffe0, this is print No.9
88 [U-MODE] pid: 3, sp is 0000003ffffffffffe0, this is print No.10
89 [U-MODE] pid: 3, sp is 0000003ffffffffffe0, this is print No.1

```

## 4.7 添加 ELF 支持

### ELF Header

ELF 文件中包含了将程序加载到内存所需的信息。当我们通过 `readElf` 来查看一个 ELF 可执行文件的时候，我们可以读到被包含在 ELF Header 中的信息：

```

1  $ readelf -a uapp
2  ELF Header:
3      Magic:   7f 45 4c 46 02 01 01 00 00 00 00 00 00 00 00
4      Class:                                ELF64
5      Data:                                    2's complement, little endian
6      Version:                                1 (current)
7      OS/ABI:                                UNIX - System V
8      ABI Version:                            0
9      Type:                                    EXEC (Executable file)
10     Machine:                                RISC-V
11     Version:                                0x1
12     Entry point address:                     0x100e8
13     Start of program headers:                 64 (bytes into file)
14     Start of section headers:                 3200 (bytes into file)
15     Flags:                                    0x0
16     Size of this header:                      64 (bytes)
17     Size of program headers:                  56 (bytes)
18     Number of program headers:                 3
19     Size of section headers:                  64 (bytes)
20     Number of section headers:                 9
21     Section header string table index:        8

```

```

22
23 Section Headers:
24   [Nr] Name                Type                Address                Offset
25         Size                EntSize                Flags Link Info Align
26   [ 0]                                NULL                0000000000000000    00000000
27         0000000000000000    0000000000000000                0    0    0
28   [ 1] .text                PROGBITS                00000000000100e8    000000e8
29         000000000000006fc    0000000000000000    AX    0    0    4
30   [ 2] .rodata                PROGBITS                00000000000107e8    000007e8
31         00000000000000078    0000000000000000    A    0    0    8
32   [ 3] .bss                  NOBITS                0000000000011860    00000860
33         000000000000003f0    0000000000000000    WA    0    0    8
34   [ 4] .comment                PROGBITS                0000000000000000    00000860
35         0000000000000002b    0000000000000001    MS    0    0    1
36   [ 5] .riscv.attributes      RISCV_ATTRIBUTE        0000000000000000    0000088b
37         0000000000000002e    0000000000000000                0    0    1
38   [ 6] .symtab                SYMTAB                0000000000000000    000008c0
39         000000000000002d0    0000000000000018                7   17    8
40   [ 7] .strtab                STRTAB                0000000000000000    00000b90
41         000000000000000a1    0000000000000000                0    0    1
42   [ 8] .shstrtab                STRTAB                0000000000000000    00000c31
43         0000000000000049    0000000000000000                0    0    1
44 Key to Flags:
45   w (write), A (alloc), X (execute), M (merge), S (strings), I (info),
46   L (link order), O (extra OS processing required), G (group), T (TLS),
47   C (compressed), x (unknown), o (OS specific), E (exclude),
48   D (mbind), p (processor specific)
49
50 There are no section groups in this file.
51
52 Program Headers:
53   Type                Offset                VirtAddr                PhysAddr
54         FileSiz                MemSiz                Flags Align
55   RISCV_ATTRIBUTE    0x000000000000088b    0x0000000000000000    0x0000000000000000
56         0x000000000000002e    0x0000000000000000    R    0x1
57   LOAD                0x00000000000000e8    0x00000000000100e8    0x00000000000100e8
58         0x0000000000000778    0x000000000001b68    RWE    0x8
59   GNU_STACK            0x0000000000000000    0x0000000000000000    0x0000000000000000
60         0x0000000000000000    0x0000000000000000    RW    0x10
61
62 Section to Segment mapping:
63   Segment Sections...
64   00      .riscv.attributes
65   01      .text .rodata .bss
66   02
67
68 There is no dynamic section in this file.
69
70 There are no relocations in this file.
71
72 The decoding of unwind sections for machine type RISC-V is not currently
73 supported.
74
75 Symbol table '.symtab' contains 30 entries:
76   Num:    Value                Size Type    Bind    vis    Ndx Name

```

```

76      0: 0000000000000000      0 NOTYPE  LOCAL  DEFAULT  UND
77      1: 000000000000100e8      0 SECTION LOCAL  DEFAULT    1 .text
78      2: 000000000000107e8      0 SECTION LOCAL  DEFAULT    2 .rodata
79      3: 00000000000011860      0 SECTION LOCAL  DEFAULT    3 .bss
80      4: 0000000000000000      0 SECTION LOCAL  DEFAULT    4 .comment
81      5: 0000000000000000      0 SECTION LOCAL  DEFAULT    5
      .riscv.attributes
82      6: 0000000000000000      0 FILE    LOCAL  DEFAULT  ABS start.o
83      7: 000000000000100e8      0 NOTYPE  LOCAL  DEFAULT    1 $x
84      8: 0000000000000000      0 FILE    LOCAL  DEFAULT  ABS getpid.c
85      9: 000000000000100ec     52 FUNC    LOCAL  DEFAULT    1 getpid
86     10: 000000000000100ec      0 NOTYPE  LOCAL  DEFAULT    1 $x
87     11: 00000000000010120      0 NOTYPE  LOCAL  DEFAULT    1 $x
88     12: 0000000000000000      0 FILE    LOCAL  DEFAULT  ABS printf.c
89     13: 0000000000001017c      0 NOTYPE  LOCAL  DEFAULT    1 $x
90     14: 000000000000101d4    1412 FUNC    LOCAL  DEFAULT    1 vprintfmt
91     15: 000000000000101d4      0 NOTYPE  LOCAL  DEFAULT    1 $x
92     16: 00000000000010758      0 NOTYPE  LOCAL  DEFAULT    1 $x
93     17: 00000000000010758     140 FUNC    GLOBAL DEFAULT    1 printf
94     18: 00000000000012060      0 NOTYPE  GLOBAL  DEFAULT  ABS
      __global_pointer$
95     19: 00000000000011860      0 NOTYPE  GLOBAL  DEFAULT    2 __SDATA_BEGIN__
96     20: 00000000000011860      4 OBJECT  GLOBAL  DEFAULT    3 tail
97     21: 000000000000100e8      0 NOTYPE  GLOBAL  DEFAULT    1 _start
98     22: 00000000000011868    1000 OBJECT  GLOBAL  DEFAULT    3 buffer
99     23: 00000000000011c50      0 NOTYPE  GLOBAL  DEFAULT    3 __BSS_END__
100    24: 00000000000011860      0 NOTYPE  GLOBAL  DEFAULT    3 __bss_start
101    25: 00000000000010120     92 FUNC    GLOBAL  DEFAULT    1 main
102    26: 0000000000001017c     88 FUNC    GLOBAL  DEFAULT    1 putc
103    27: 00000000000011860      0 NOTYPE  GLOBAL  DEFAULT    2 __DATA_BEGIN__
104    28: 00000000000011860      0 NOTYPE  GLOBAL  DEFAULT    2 _edata
105    29: 00000000000011c50      0 NOTYPE  GLOBAL  DEFAULT    3 _end
106
107    No version information found in this file.
108    Attribute Section: riscv
109    File Attributes
110      Tag_RISCV_arch: "rv64i2p0_m2p0_a2p0_f2p0_d2p0"

```

其中包含了两种将程序分块的粒度，Segment（段）和 Section（节），我们以段为粒度将程序加载进内存中。可以看到，给出的样例程序包含了三个段，这里我们只关注 Type 为 LOAD 的段，LOAD 表示他们需要在开始运行前被加载进内存中，这是我们在初始化进程的时候需要执行的工作。

而“节”代表了更细分的语义，比如 `.text` 一般包含了程序的指令，`.rodata` 是只读的全局变量等，大家可以自行 Google 来学习更多相关内容。

## 所以代码怎么写？

首先我们需要将 `uapp.S` 中的 payload 给换成我们的 ELF 文件。

```

1  /* user/uapp.S */
2  .section .uapp
3
4  .incbin "uapp"
5

```

这时候从 `uapp_start` 开始的数据就变成了名为 `uapp` 的 ELF 文件，也就是说 `uapp_start` 处 32-bit 的数据不再是我们需要执行第一条指令了，而是 ELF Header 的开始。

这时候就需要你对 `task_init` 中的初始化步骤进行修改。我们给出了 ELF 相关的结构体定义 (`elf.h`)，大家可以直接使用。你可能会使用到的结构体或者域如下：

```
1  Elf64_Ehdr    // 你可以将 uapp_start 强制转化为改类型的指针，
2                然后把那一块内存当成此类结构体来读其中的数据，其中包括：
3      e_ident    // Magic Number，你可以通过这个域来检测自己是不是真的正在读一个 Ehdr，
4                值一定是 7f 45 4c 46 02 01 01 00 00 00 00 00 00 00 00
5      e_entry     // 程序的第一条指令被存储的用户态虚拟地址
6      e_phnum     // ELF 文件包含的 Segment 的数量
7      e_phoff     // ELF 文件包含的 Segment 数组相对于 Ehdr 的偏移量
8
9  Elf64_Phdr    // 存储了程序各个 Segment 相关的 metadata
10              // 你可以将 uapp_start + e_phoff 强制转化为此类型，就会指向第一个
    Phdr，
11              // uapp_start + e_phoff + 1 * sizeof(Elf64_Phdr)，则指向第二个
12      p_filesz   // Segment 在文件中占的大小
13      p_memsz    // Segment 在内存中占的大小
14      p_vaddr    // Segment 起始的用户态虚拟地址
15      p_offset   // Segment 在文件中相对于 Ehdr 的偏移量
16      p_type     // Segment 的类型
17      p_flags    // Segment 的权限（包括了读、写和执行）
18
```

我们可以按照这些信息，在从 `uapp_start` - `uapp_end` 这个 ELF 文件中的内容 **拷贝** 到我们开辟的内存中。

其中相对文件偏移 `p_offset` 指出相应 segment 的内容从 ELF 文件的第 `p_offset` 字节开始，在文件中的大小为 `p_filesz`，它需要被分配到以 `p_vaddr` 为首地址的虚拟内存位置，在内存中它占用大小为 `p_memsz`。也就是说，这个 segment 使用的内存就是 `[p_vaddr, p_vaddr + p_memsz)` 这一连续区间，然后将 segment 的内容从 ELF 文件中读入到这一内存区间，并将 `[p_vaddr + p_filesz, p_vaddr + p_memsz)` 对应的物理区间清零。（本段内容引用自[南京大学PA](#)）

你也可以参考这篇 [blog](#) 中关于 **静态** 链接程序的载入过程来进行你的载入。

这里有不少例子可以举，为了避免同学们在实验中花太多时间，我们告诉大家可以怎么找到实验中这些相关变量被存在了哪里：（注意以下的 `uapp_start` 类型使用的是 `char*`，如果你在使用其他类型，需要根据你使用的类型去调整针对指针的算数运算。）

- `Elf64_Ehdr* ehdr = (Elf64_Ehdr*)uapp_start`，从地址 `uapp_start` 开始，便是我们要找的 Ehdr。
- `Elf64_Phdr* phdrs = (Elf64_Phdr*)(uapp_start + ehdr->phoff)`，是一个 Phdr 数组，其中的每个元素都是一个 `Elf64_Phdr`。
- `phdrs[ehdr->phnum - 1]` 是最后一个 Phdr。
- `phdrs[0].p_type == PT_LOAD`，说明这个 Segment 的类型是 LOAD，需要在初始化时被加载进内存。
- `(void*)(uapp_start + phdrs[1].p_offset)` 将会指向第二个段中的内容的开头。
- `load_program` 具体实现如下：
  - 在实现过程中由于每个 segment 将要分配到的虚拟地址没有与页面对齐，所以需要进行对齐

- o 在实现中发现, 需要将对齐时多余的内存也要进行分配才能完成所有 segment 的对齐

```
1 // /arch/riscv/kernel/proc.c
2 static uint64_t load_program(struct task_struct* task) {
3     Elf64_Ehdr* ehdr = (Elf64_Ehdr*)uapp_start;
4
5     uint64_t phdr_start = (uint64_t)ehdr + ehdr->e_phoff;
6     int phdr_cnt = ehdr->e_phnum;
7
8     Elf64_Phdr* phdr;
9     int load_phdr_cnt = 0;
10    for (int i = 0; i < phdr_cnt; i++) {
11        phdr = (Elf64_Phdr*)(phdr_start + sizeof(Elf64_Phdr) * i);
12        if (phdr->p_type == PT_LOAD) {
13            // alloc space and copy content
14            uint64 page_num = (phdr->p_vaddr - PGROUNDDOWN((uint64)phdr->p_vaddr) + phdr->p_memsz + PGSIZE - 1) / PGSIZE;
15            uint64 uapp = alloc_pages(page_num);
16            memset(uapp, 0, page_num*PGSIZE);
17            for(int t = 0; t < page_num; t++){
18                uint64* src = (uint64*)(uapp_start + t * PGSIZE + phdr->p_offset);
19                uint64* dst = (uint64*)(uapp + t * PGSIZE + phdr->p_vaddr - PGROUNDDOWN((uint64)phdr->p_vaddr));
20                for (int k = 0; k < PGSIZE / 8; k++) {
21                    dst[k] = src[k];
22                }
23            }
24            // do mapping X|W|R|V
25            create_mapping(task->pgd, PGROUNDDOWN((uint64)phdr->p_vaddr), uapp - PA2VA_OFFSET, PGSIZE*page_num,
26                ((phdr->p_flags & PF_X) << 3) | ((phdr->p_flags & PF_W) << 1) | ((phdr->p_flags & PF_R) >> 1) | 0b10001);
27        }
28    }
29
30    // allocate user stack and do mapping
31    // code...
32    create_mapping(task->pgd, USER_END - PGSIZE, (uint64)kalloc - PA2VA_OFFSET, PGSIZE, 23);
33
34    // following code has been written for you
35    // set user stack
36    // pc for the user program
37    task->thread.sepc = ehdr->e_entry;
38    // sstatus bits set
39    task->thread.sstatus = (1ULL << 5) + (1ULL << 18);
40    // user stack for user program
41    task->thread.sscratch = USER_END;
42 }
```

## 4. 思考题

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）

用户态（一）对内核态（一）

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

`p_filesz` 通常表示段在文件中的大小，而 `p_memsz` 表示段在内存中的大小。大小不同的原因可能有如下三点：

- 在执行加载程序时，OS可能要求段在内存中的大小是某个对齐值的倍数，就会导致 `p_memsz` 大于 `p_filesz`，多分配的空间会被填充为0
- BSS段的全局变量在可执行文件中不占用实际空间，仅仅记录大小，而在内存中则需要分配空间并初始化。
- 有些段可能存在冬天分配生成的数据，这些数据在文件中是不存在的

3. 为什么多个进程的栈虚拟地址可以是相同的？用户有没有常规的方法知道自己栈所在的物理地址？

因为每个进程都有自己的专属页表，在不同的进程页表中，进程的栈的虚拟地址映射到的物理地址是不相同的。

在Linux系统中，`/proc/pid/pagemap` 文件记录了对应进程的物理页号信息，我们可以通过该文件查看每个虚拟页映射到的物理页

```
1 # proc/pid/pagemap. This file lets a userspace process find out which
2   physical frame each virtual page is mapped to. It contains one 64-
   bit
3   value for each virtual page, containing the following data (from
4   fs/proc/task_mmu.c, above pagemap_read):
5
6   * Bits 0-54  page frame number (PFN) if present//present为1时, bit0-
   54表示物理页号
7   * Bits 0-4   swap type if swapped
8   * Bits 5-54  swap offset if swapped
9   * Bit 55     pte is soft-dirty (see Documentation/vm/soft-dirty.txt)
10  * Bit 56     page exclusively mapped (since 4.2)
11  * Bits 57-60 zero
12  * Bit 61     page is file-page or shared-anon (since 3.5)
13  * Bit 62     page swapped
14  * Bit 63     page present//如果为1, 表示当前物理页在内存中; 为0, 表示当前物
   理页不在内存中
```

## 5. 心得体会

在本次做实验的过程中，由于要修改的文件较多，所以才debug过程中十分艰难，在进trap的时候，一直报scause=12的错误，即Instruction Fetch Fault。于是去Debug初始化的部分和汇编文件，但是均发现没有什么奇怪的地方(表面上)，后来逐步Debug发现每次跳到用户空间第一个地址就寄了，于是就去查看页表的创建。发现在具体的创建过程中uapp.bin文件根本没有加载到内存中去..... 最后终于发现pull下来的文件和网站上不一样。

具体实践的过程中，发现自己对于步骤的理解并不十分清楚，在实现过程中不甚顺利，还得加强自己的编程能力。



