

Lab 4: RV64 用户态程序

姓名：高铭健

学号：3210102322

日期：2023.12.10

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

- 计算机（Intel Core i5以上，4GB内存以上）系统
- Ubuntu 22.04.2 LTS

3. 背景知识

3.1 用户模式和内核模式

处理器存在两种不同的模式：**用户模式**（U-Mode）和**内核模式**（S-Mode）。

- 在用户模式下，执行代码无法直接访问硬件，必须委托给系统提供的接口才能访问硬件或内存。
- 在内核模式下，执行代码对底层硬件具有完整且不受限制的访问权限，它可以执行任何 CPU 指令并引用任何内存地址。

处理器根据处理器上运行的代码类型在这两种模式之间切换。应用程序以用户模式运行，而核心操作系统组件以内核模式运行。

3.2 目标

- 当启动用户态应用程序时，内核将为该应用程序创建一个进程，并提供了专用虚拟地址空间等资源。
 - 每个应用程序的虚拟地址空间是私有的，一个应用程序无法更改属于另一个应用程序的数据。
 - 每个应用程序都是独立运行的，如果一个应用程序崩溃，其他应用程序和操作系统将不会受到影响。
- 用户态应用程序可访问的虚拟地址空间是受限的。
 - 在用户态下，应用程序无法访问内核的虚拟地址，防止其修改关键操作系统数据。
 - 当用户态程序需要访问关键资源的时候，可以通过**系统调用**来完成用户态程序与操作系统之间的互动。

3.3 系统调用约定

系统调用是用户态应用程序请求内核服务的一种方式。在 RISC-V 中，我们使用 `ecall` 指令进行系统调用。当执行这条指令时，处理器会提升特权模式，跳转到异常处理函数，处理这条系统调用。

3.4 sstatus[SUM] PTE[U]

当页表项 PTE[U] 置 0 时，该页表项对应的内存页为内核页，运行在 U-Mode 下的代码无法访问。当页表项 PTE[U] 置 1 时，该页表项对应的内存页为用户页，运行在 S-Mode 下的代码无法访问。如果想让 S 特权级下的程序能够访问用户页，需要对 sstatus[SUM] 位置 1。但是无论什么样的情况下，用户页中的指令对于 S-Mode 而言都是**无法执行的**。

3.5 用户态栈与内核态栈

当用户态程序在用户态运行时，其使用的栈为**用户态栈**。当进行系统调用时，陷入内核处理时使用的栈为**内核态栈**。因此，需要区分用户态栈和内核态栈，且需要在异常处理的过程中对栈进行切换。

3.6 ELF 程序

ELF, short for Executable and Linkable Format. 是当今被广泛使用的应用程序格式。例如当我们运行 `gcc <some-name>.c` 后产生的 `a.out` 输出文件的格式就是 ELF。

```
1  $ cat hello.c
2  #include <stdio.h>
3
4  int main() {
5      printf("hello, world\n");
6      return 0;
7  }
8  $ gcc hello.c
9  $ file a.out
10 a.out: ELF 64-bit LSB pie executable, x86-64, version 1 (SYSV), dynamically
    linked,
11 interpreter /lib64/ld-linux-x86-64.so.2,
    BuildID[sha1]=dd33139196142abd22542134c20d85c571a78b0c,
12 for GNU/Linux 3.2.0, not stripped
```

将程序封装成 ELF 格式的意义包括以下几点：

- ELF 文件可以包含将程序正确加载入内存的元数据（metadata）。
- ELF 文件在运行时可以由加载器（loader）将动态链接在程序上的动态链接库（shared library）正确地
从硬盘或内存中加载。
- ELF 文件包含的重定位信息可以让该程序继续和其他可重定位文件和库再次链接，构成新的可执行文件。

4. 实验步骤

4.1 准备工程

- 修改 `vmlinux.lds`, 将用户态程序 `uapp` 加载至 `.data` 段。添加如下片段:

```
1  .data : ALIGN(0x1000){
2      _sdata = .;
3
4      *(.sdata .sdata*)
5      *(.data .data.*)
6
7      _edata = .;
8
9      . = ALIGN(0x1000);
10     _sramdisk = .;
11     *(.uapp .uapp*)
12     _eramdisk = .;
13     . = ALIGN(0x1000);
14
15 }
```

- 修改 `defs.h`, 在 `defs.h` 添加 USER space的起始虚拟地址和结尾虚拟地址:

```
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:

```
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, 添加 `${MAKE} -C user all`; 代码如下所示:

```

1  export
2  CROSS_=riscv64-linux-gnu-
3  GCC=${CROSS_}gcc
4  LD=${CROSS_}ld
5  OBJCOPY=${CROSS_}objcopy
6
7  ISA=rv64imafdzifencei
8  ABI=lp64
9
10 INCLUDE = -I $(shell pwd)/include -I $(shell pwd)/arch/riscv/include
11 CF = -march=$(ISA) -mabi=$(ABI) -mcmode=medany -fno-builtin -ffunction-sections
    -fdata-sections -nostartfiles -nostdlib -nostdinc -static -lgcc -wl,--nmagic -
    wl,--gc-sections -g
12 CFLAG = ${CF} ${INCLUDE}
13
14 .PHONY:all run debug clean
15 all: clean
16     ${MAKE} -C lib all
17     ${MAKE} -C test all
18     ${MAKE} -C init all
19     ${MAKE} -C arch/riscv all
20     ${MAKE} -C user all
21     @echo -e '\n'Build Finished OK
22
23 TEST:
24     ${MAKE} -C lib all
25     ${MAKE} -C test test
26     ${MAKE} -C init all
27     ${MAKE} -C user all
28     ${MAKE} -C arch/riscv all
29     @echo -e '\n'Build Finished OK
30
31 run: all
32     @echo Launch the qemu .....
33     @qemu-system-riscv64 -nographic -machine virt -kernel vmlinux -bios default
34
35 test-run: TEST
36     @echo Launch the qemu .....
37     @qemu-system-riscv64 -nographic -machine virt -kernel vmlinux -bios default
38
39 debug: all
40     @echo Launch the qemu for debug .....
41     @qemu-system-riscv64 -nographic -machine virt -kernel vmlinux -bios default
    -S -s
42
43 test-debug: TEST
44     @echo Launch the qemu for debug .....
45     @qemu-system-riscv64 -nographic -machine virt -kernel vmlinux -bios default
    -S -s
46
47 clean:
48     ${MAKE} -C lib clean

```

```

49     ${MAKE} -C test clean
50     ${MAKE} -C init clean
51     ${MAKE} -C arch/riscv clean
52     $(shell test -f vmlinux && rm vmlinux)
53     $(shell test -f System.map && rm System.map)
54     @echo -e '\n'Clean Finished

```

4.2 创建用户态进程

- 创建 4 个用户态进程，修改 `proc.h` 中的 `NR_TASKS`：

```

1  #define NR_TASKS (1 + 3)

```

- 用户态进程要对 `sepc` `sstatus` `sscratch` 做设置，将其加入 `thread_struct` 中，同时将内核和用户栈指针以及 `satp` 都直接添加到 `task_struct` 方便处理。因为每一个进程都创建一个页表，所以添加 `pagetable_t pgd`（没有使用 `struct thread_info *thread_info`，使之前的 `switch_to` 中变量的偏移量不用修改）：

```

1
2  struct thread_struct {
3      uint64 ra;
4      uint64 sp;
5      uint64 s[12];
6
7      uint64 sepc, sstatus, sscratch;
8  };
9
10 struct task_struct {
11     struct thread_info *thread_info;
12     uint64 state;    // 线程状态
13     uint64 counter;  // 运行剩余时间
14     uint64 priority; // 运行优先级 1最低 10最高
15     uint64 pid;      // 线程id
16     struct thread_struct thread;
17     uint64 satp;
18     uint64 kernel_sp;
19     uint64 user_sp;
20     pagetable_t pgd;
21 };

```

- 修改 `task_init`
 - 对于每个进程，初始化在 `thread_struct` 中添加的三个变量：
 - 将 `sepc` 设置为 `USER_START`。
 - 将 `sstatus` 中的 `SPP` 置0（使得 `sret` 返回至 U-Mode），`SPIE` 置1（`sret` 之后开启中断），`SUM` 置1（S-Mode 可以访问 User 页面）。

XLEN-1		XLEN-2										20	19	18	17
SD		0										MXR		SUM	0
1		XLEN-21										1		1	1
16	15	14	13	12	9	8	7	6	5	4	3	2	1	0	
XS[1:0]		FS[1:0]		0		SPP	0		SPIE	UPIE	0		SIE	UIE	
2		2		4		1	2		1	1	2		1	1	

- 将 `sscratch` 设置为 U-Mode 的 `sp`，其值为 `USER_END`（即，用户态栈被放置在 `user space` 的最后一个页面）。
- 直接设置 `satp`，将 `PPN` 清空，并设置为当前进程申请的页表的物理页号

```

1 //将 sepc 设置为 USER_START
2 task[i]->thread.sepc = USER_START;
3 //配置 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）
4 uint64 sstatus = csr_read(sstatus);
5 sstatus |= 0x0000000000040020;
6 sstatus &= 0xffffffffffffefff;
7 task[i]->thread.sstatus = sstatus;
8 //将 sscratch 设置为 U-Mode 的 sp，其值为 USER_END（即，用户态栈被放置在 user
  space 的最后一个页面）
9 task[i]->thread.sscratch = USER_END;
10 task[i]->satp = (csr_read(satp) & 0xfffff00000000000) |
  (((uint64)new_pgtbl - PA2VA_OFFSET) >> 12);

```

- 对于每个进程，创建属于它自己的页表，使用 `alloc_page()` 函数进行内存申请：

```

1 //对于每个进程，创建属于它自己的页表
2 uint64* new_pgtbl = (uint64*)alloc_page();

```

- 为了避免 U-Mode 和 S-Mode 切换的时候切换页表，将内核页表（`swapper_pg_dir`）复制到每个进程的页表中，注意在前面声明的 `_sramdisk[]` 类型应该是 `uint64`，然后使用for循环复制一个页的大小即可。

```

1 extern uint64 _sramdisk[];
2 ...
3 for(unsigned long i = 0; i < PGSIZE; i++){
4     new_pgtbl[i] = swapper_pg_dir[i];
5 }

```

- 将 `uapp` 所在的页面映射到每个进行的页表中，在程序运行过程中，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。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，造成预期外的进程间相互影响。

使用 `alloc_page()` 函数申请一片内存，将 `_sramdisk` 之后的文件复制到这片内存中

然后调用 `create_mapping()` 函数进行映射，将 `USER_START` 开始的user指令映射到申请的内存中，大小为1个页，权限是V=1（有效） | R=1（可读） | W=1（可写） | X=1（可执行） | U=1（说明是用户页）

```
1 //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进行映射，防止所有的进程共享数据
2 uint64* new_addrress = (uint64*)alloc_page();
3 for(unsigned long i = 0;i < PGSIZE;i++){
4     new_addrress[i] = _sramdisk[i];
5 }
6 //将 uapp 所在的页面映射到每个进行的页表中
7 uint64 va = USER_START;
8 uint64 pa = (uint64)new_addrress - PA2VA_OFFSET;
9 uint64 sz = PGSIZE;
10 int perm = 0b11111;
11 printk("task[i]->uapp %|x %|x %|x %d\n",va,pa,sz,perm);
12 create_mapping(new_pgtbl, va, pa, sz, perm);
```

- 设置用户态栈。对每个用户态进程，其拥有两个栈：用户态栈和内核态栈；通过 `alloc_page` 接口申请一个空的页面来作为用户态栈，然后调用 `create_mapping()` 函数进行映射，将 `USER_END - PGSIZE` 开始的use页映射到申请的内存中，权限是V=1（有效） | R=1（可读） | W=1（可写） | X=0（用户页中的指令对于 S-Mode 而言是**无法执行的**） | U=1（说明是用户页）

```
1 //设置用户态栈,通过 alloc_page 接口申请一个空的页面来作为用户态栈，并映射到进程的页表中
2 task[i]->user_sp = alloc_page();
3 va = USER_END - PGSIZE;
4 pa = task[i]->user_sp - PA2VA_OFFSET;
5 sz = PGSIZE;
6 perm = 0b10111;
7 printk("task[i]->user %|x %|x %|x %d\n",va,pa,sz,perm);
8 create_mapping(new_pgtbl, va, pa, sz, perm);
```

- 修改 `__switch_to`，在切换进程时添加保存进程的 `sepc sstatus sscratch satp`，同样LD时获得新进程的这些值

```
1 __switch_to:
2     #加入保存/恢复 sepc sstatus sscratch 以及 切换页表的逻辑
3     csrr t1, sepc
4     csrr t2, sstatus
5     csrr t3, sscratch
6     csrr t4, satp
7     # save state to prev process
8     # YOUR CODE HERE
9     addi t0, a0, 40
10    sd ra, 0(t0)
11    sd sp, 8(t0)
```

```

12     sd s0, 16(t0)
13     sd s1, 24(t0)
14     sd s2, 32(t0)
15     sd s3, 40(t0)
16     sd s4, 48(t0)
17     sd s5, 56(t0)
18     sd s6, 64(t0)
19     sd s7, 72(t0)
20     sd s8, 80(t0)
21     sd s9, 88(t0)
22     sd s10, 96(t0)
23     sd s11, 104(t0)
24     sd t1, 112(t0)
25     sd t2, 120(t0)
26     sd t3, 128(t0)
27     sd t4, 136(t0)
28     # restore state from next process
29     # YOUR CODE HERE
30     addi t0, a1, 40
31     ld ra, 0(t0)
32     ld sp, 8(t0)
33     ld s0, 16(t0)
34     ld s1, 24(t0)
35     ld s2, 32(t0)
36     ld s3, 40(t0)
37     ld s4, 48(t0)
38     ld s5, 56(t0)
39     ld s6, 64(t0)
40     ld s7, 72(t0)
41     ld s8, 80(t0)
42     ld s9, 88(t0)
43     ld s10, 96(t0)
44     ld s11, 104(t0)
45     ld t1, 112(t0)
46     ld t2, 120(t0)
47     ld t3, 128(t0)
48     ld t4, 136(t0)
49     csrw sepc, t1
50     csrw sstatus, t2
51     csrw sscratch, t3
52     csrw satp, t4
53     # flush
54     sfence.vma zero, zero
55     ret

```

- 在切换页表之后，通过 `fence.i` 和 `vma.fence` 来刷新 TLB 和 ICache：


```

1  ...
2  # flush tlb
3  sfence.vma zero, zero
4  # flush icache
5  fence.i

```

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

- 修改 __dummy, 交换 thread_struct.sp 和保存在 thread_struct.sscratch 的 U-Mode sp 来完成进程栈的切换, 同时将 sepc 设为0, 也就是要返回 user space 的开头执行指令:

```

1  __dummy:
2      # YOUR CODE HERE
3      #交换对应的寄存器的值
4      csrr t1, sscratch
5      mv t2, sp
6      mv sp, t1
7      csrw sscratch, t2
8      # sepc置0
9      addi t0, x0, 0
10     csrw sepc, t0
11     sret #从中断中返回

```

- 修改 _trap, 在 _trap 的首尾进行内核和用户的切换, 但是如果已经在内核态那么显然不用切换, 判断当前的 sscratch 是否为 0, 是的话就说明已经在内核态, 不用切换

```

1  _traps:
2      # YOUR CODE HERE
3      # -----
4      #在 _trap 的首尾我们都需要做类似的操作(判断是否在内核态)
5      csrr t0, sscratch
6      addi t1, x0, 0
7      beq t0, t1, Return
8      #交换对应的寄存器的值
9      csrr t1, sscratch
10     mv t2, sp
11     mv sp, t1
12     csrw sscratch, t2
13     # 1. save 32 registers and sepc to stack
14     Next:
15     addi sp, sp, -280
16     #存入32个寄存器
17     sd x0, 0(sp)
18     sd x1, 8(sp)
19     sd x2, 16(sp)

```

```

20      sd x3, 24(sp)
21      sd x4, 32(sp)
22      sd x5, 40(sp)
23      sd x6, 48(sp)
24      sd x7, 56(sp)
25      sd x8, 64(sp)
26      sd x9, 72(sp)
27      sd x10, 80(sp) #寄存器a0
28      sd x11, 88(sp) #寄存器a1
29      sd x12, 96(sp)
30      sd x13, 104(sp)
31      sd x14, 112(sp)
32      sd x15, 120(sp)
33      sd x16, 128(sp)
34      sd x17, 136(sp)
35      sd x18, 144(sp)
36      sd x19, 152(sp)
37      sd x20, 160(sp)
38      sd x21, 168(sp)
39      sd x22, 176(sp)
40      sd x23, 184(sp)
41      sd x24, 192(sp)
42      sd x25, 200(sp)
43      sd x26, 208(sp)
44      sd x27, 216(sp)
45      sd x28, 224(sp)
46      sd x29, 232(sp)
47      sd x30, 240(sp)
48      sd x31, 248(sp)
49      #存sepc
50      csrr t0, sepc
51      sd t0, 256(sp)
52      csrr t0, sstatus
53      sd t0, 264(sp)
54      csrr t0, sscratch
55      sd t0, 272(sp)
56
57
58      # -----
59
60      # 2. call trap_handler
61      csrr a0, scause
62      csrr a1, sepc
63      mv a2, sp #寄存器a2
64      jal trap_handler
65
66      # -----
67
68      # 3. restore sepc and 32 registers (x2(sp) should be restore last)
        from stack
69
70      ld t0, 272(sp)

```

```

71     csrw sscratch, t0
72     ld t0, 264(sp)
73     csrw sstatus, t0
74     ld t0, 256(sp)
75     csrw sepc, t0
76     ld x0, 0(sp)
77     ld x1, 8(sp)
78     ld x2, 16(sp)
79     ld x3, 24(sp)
80     ld x4, 32(sp)
81     ld x5, 40(sp)
82     ld x6, 48(sp)
83     ld x7, 56(sp)
84     ld x8, 64(sp)
85     ld x9, 72(sp)
86     ld x10, 80(sp) #寄存器a0
87     ld x11, 88(sp) #寄存器a1
88     ld x12, 96(sp)
89     ld x13, 104(sp)
90     ld x14, 112(sp)
91     ld x15, 120(sp)
92     ld x16, 128(sp)
93     ld x17, 136(sp)
94     ld x18, 144(sp)
95     ld x19, 152(sp)
96     ld x20, 160(sp)
97     ld x21, 168(sp)
98     ld x22, 176(sp)
99     ld x23, 184(sp)
100    ld x24, 192(sp)
101    ld x25, 200(sp)
102    ld x26, 208(sp)
103    ld x27, 216(sp)
104    ld x28, 224(sp)
105    ld x29, 232(sp)
106    ld x30, 240(sp)
107    ld x31, 248(sp)
108
109    addi sp, sp, 280
110    # -----
111    #在 _trap 的首尾我们都需要做类似的操作(判断是否在内核态)
112    csrr t0, sscratch
113    addi t1, x0, 0
114    beq t0, t1, Return
115    #交换对应的寄存器的值
116    csrr t1, sscratch
117    mv t2, sp
118    mv sp, t1
119    csrw sscratch, t2
120    Return:
121    # 4. return from trap
122    sret

```

- 在 `trap_handler` 里面进行捕获用户态的 `ecall`。修改 `trap_handler`，传入 `pt_regs`，其中包括了全部的寄存器，同时定义需要 `SYS_WRITE` 和 `SYS_PID` 的对应寄存器的值。

首先根据 `scause` 的值是否为 8 判断是否进行 `ecall` 处理，然后根据 `a7` 寄存器的值判断需要 `SYS_WRITE` 还是 `SYS_PID`，如果是 64 则调用 `sys_write()` 函数，传入需要输出的字符串 `a1` 和大小 `a2`；如果是 172 则调用 `sys_getpid()` 函数，获得当前 `pid` 并赋给 `a0`（修改 `regs->reg[10]`）

最后将 `sepc + 4` 返回异常指令的下一条指令

```
1  #define SYS_WRITE 64
2  #define SYS_PID 172
3  struct pt_regs
4  {
5      uint64 reg[32];
6      uint64 sepc;
7      uint64 sstatus;
8      uint64 sscratch;
9  };
10 void trap_handler(unsigned long scause, unsigned long sepc, struct pt_regs
    *regs) {
11     // 通过 `scause` 判断trap类型
12     // 如果是interrupt 判断是否是timer interrupt
13     // 如果是timer interrupt 则打印输出相关信息，并通过 `clock_set_next_event()`
    设置下一次时钟中断
14     if(scause & 0x8000000000000000){//trap 类型为interrupt
15         if(scause == 0x8000000000000005){//timer interrupt
16             printk("[S] Supervisor Mode Timer Interrupt\n");
17             clock_set_next_event();
18             do_timer();
19         }
20     }
21     else if(scause == 8){
22         if(regs->reg[17] == SYS_WRITE){
23             char *buffer = (char *) (regs->reg[11]);
24             int cnt = sys_write(1, buffer, regs->reg[12]);
25             regs->sepc += 4;
26         }
27         else if(regs->reg[17] == SYS_PID){
28             uint64 pid = sys_getpid();
29             regs->sepc += 4;
30             regs->reg[10] = pid;
31             //         printk("pid: %d", pid);
32         }
33         else{
34             //         printk("error!\n");
35         }
36     }
37 }
```

4.4 添加系统调用

- 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。
- 增加 `syscall.c` `syscall.h` 文件，`sys_write()` 简单输出传入的字符串，`sys_getpid()` 需要先使用当前进程 `task_struct *current` 直接返回其pid

```
1 //
2 // syscall.h
3 //
4 #include "proc.h"
5 #include "sbi.h"
6 #include "printk.h"
7
8
9 int sys_write(unsigned int fd, const char* buf, size_t count);
10
11 unsigned long sys_getpid();
```

```
1 //
2 // syscall.c
3 //
4 #include "syscall.h"
5 extern struct task_struct *current;
6 int sys_write(unsigned int fd, const char *buf, size_t count){
7     int write_size;
8     for(int i = 0; i < count; i++){
9         printk("%c", buf[i]);
10        write_size++;
11    }
12    return write_size;
13 }
14
15 unsigned long sys_getpid(){
16     return (unsigned long)(current->pid);
17 }
```

4.5 修改 head.S 以及 start_kernel

- 在 start_kernel 中调用 schedule()，放在 test() 之前：

```

1  #include "printk.h"
2  #include "sbi.h"
3  extern void test();
4  extern void schedule();
5  int start_kernel() {
6      printk("2023");
7      printk(" Hello RISC-V\n");
8      schedule();
9      test(); // DO NOT DELETE !!!
10     return 0;
11 }
12

```

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

```

1      # set sstatus[SIE] = 1
2      #addi t0, x0, 1
3      #slli t0, t0, 1
4      #csrr t1, sstatus
5      #or t2, t0, t1
6      #csrw sstatus, t2

```

4.6 测试

- 输出结果如下

```

sept= 00010000...proc_init done:
2023 Hello RISC-V

switch to [PID = 4 COUNTER = 1 PRIORITY = 66]
[U-MODE] pid: 4, sp is 0000003ffffffe0, this is print No.1

switch to [PID = 3 COUNTER = 1 PRIORITY = 52]
[U-MODE] pid: 3, sp is 0000003ffffffe0, this is print No.1

switch to [PID = 2 COUNTER = 1 PRIORITY = 88]
[U-MODE] pid: 2, sp is 0000003ffffffe0, this is print No.1

switch to [PID = 1 COUNTER = 1 PRIORITY = 37]
[U-MODE] pid: 1, sp is 0000003ffffffe0, this is print No.1

switch to [PID = 2 COUNTER = 88 PRIORITY = 88]
[U-MODE] pid: 2, sp is 0000003ffffffe0, this is print No.2
[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
QEMU: Terminated
fengnian@MICROS0-VU1H310:/mnt/d/study/OS/os23fall-stu/src/lab4$

```

这里将第一轮调度的时间片都设为1，打印结束后进入第二轮，根据打印出的No.*可以看出各进程是分离的

4.7 添加 ELF 支持

- 首先我们需要将 `uapp.S` 中的 payload 给换成 ELF 文件：

```
1  .section .uapp
2
3  .incbin "uapp"
4
```

- 将程序 load 进入内存。

先将内核页表复制到每个进程的页表中（与之前一样）

将 segment 的内容从ELF文件中读入 `[p_vaddr, p_vaddr + p_memsz)` 内存区间，先调用 `alloc_pages(page_cnt)` 申请足够页数的内存空间，然后将 `_sramdisk - _eramdisk` 这一段文件复制到刚才申请的内存（注意这里使用的是 `char _sramdisk[]`），之后调用 `create_mapping()` 进行映射，权限与之前相同

用户态栈的设置和之前相同

```
1  extern char _sramdisk[];
2  static uint64_t load_program(struct task_struct* task) {
3      Elf64_Ehdr* ehdr = (Elf64_Ehdr*)_sramdisk;
4
5      uint64_t phdr_start = (uint64_t)ehdr + ehdr->e_phoff;
6      int phdr_cnt = ehdr->e_phnum;
7      uint64* new_pgtbl = (uint64*)alloc_page();
8      //将内核页表 ( swapper_pg_dir ) 复制到每个进程的页表中
9      for(unsigned long i = 0; i < PGSIZE; i++){
10         new_pgtbl[i] = swapper_pg_dir[i];
11     }
12
13     Elf64_Phdr* phdr;
14     int load_phdr_cnt = 0;
15     for (int i = 0; i < phdr_cnt; i++) {
16         phdr = (Elf64_Phdr*)(phdr_start + sizeof(Elf64_Phdr) * i);
17         if (phdr->p_type == PT_LOAD) {
18             // alloc space and copy content
19             int page_cnt;
20             if(phdr->p_memsz % PGSIZE == 0){
21                 page_cnt = phdr->p_memsz / PGSIZE;
22             }
23             else{
24                 page_cnt = phdr->p_memsz / PGSIZE + 1;
25             }
26             uint64* new_space = (uint64 *)alloc_pages(page_cnt);
27             uint64* src = (uint64 *)(_sramdisk);
28             memcpy(new_space, src, phdr->p_memsz);
29             // do mapping
```

```

30         uint64 va = phdr->p_vaddr;
31         uint64 pa = (uint64)new_space - PA2VA_OFFSET;
32         create_mapping(new_pgtbl, va, pa, phdr->p_memsz, 31);
33     }
34 }
35 //设置用户态栈,通过 alloc_page 接口申请一个空的页面来作为用户态栈,并映射到进程的页
表中
36 task->user_sp = alloc_page();
37 uint64 va = USER_END - PGSIZE;
38 uint64 pa = task->user_sp - PA2VA_OFFSET;
39 uint64 sz = PGSIZE;
40 uint64 perm = 0b10111;
41 create_mapping(new_pgtbl, va, pa, sz, perm);
42 task->satp = (csr_read(satp) & 0xfffff00000000000) |
(((uint64)new_pgtbl - PA2VA_OFFSET) >> 12);
43 //将 sepc 设置为 ehdr->e_entry
44 task->thread.sepc = ehdr->e_entry;
45 printf("sepc= %lx", ehdr->e_entry);
46 //配置 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)
47 uint64 sstatus = csr_read(sstatus);
48 sstatus |= 0x0000000000040020;
49 sstatus &= 0xffffffffffffeff;
50 task->thread.sstatus = sstatus;
51 //将 sscratch 设置为 U-Mode 的 sp, 其值为 USER_END (即, 用户态栈被放置在 user
space 的最后一个页面)
52 task->thread.sscratch = USER_END;
53 }

```

- 将之前的 task_init 替换, 直接调用 load_program(task[i]):

```

1 void task_init() {
2     test_init(NR_TASKS);
3     // 1. 调用 kalloc() 为 idle 分配一个物理页
4     idle = kalloc();
5     // 2. 设置 state 为 TASK_RUNNING;
6     idle->state = TASK_RUNNING;
7     // 3. 由于 idle 不参与调度 可以将其 counter / priority 设置为 0
8     idle->counter = 0;
9     idle->priority = 0;
10    // 4. 设置 idle 的 pid 为 0
11    idle->pid = 0;
12    // 5. 将 current 和 task[0] 指向 idle
13    current = idle;
14    task[0] = idle;
15    for(int i = 1; i < NR_TASKS; i++){
16        task[i] = kalloc();
17        task[i]->pid = i;
18        task[i]->counter = task_test_counter[i];
19        task[i]->priority = task_test_priority[i];
20        task[i]->state = TASK_RUNNING;

```



```

21     uint64 ra_address = &__dummy;
22     uint64 sp_address = (uint64)task[i] + PGSIZE;
23     task[i]->thread.ra = ra_address;
24     task[i]->thread.sp = sp_address;
25     task[i]->kernel_sp = sp_address;
26     load_program(task[i]);
27 }
28 }
29

```

- 添加 ELF 支持后的测试结果：

```

2023 Hello RISC-V

switch to [PID = 4 COUNTER = 1 PRIORITY = 66]
[U-MODE] pid: 4, sp is 0000003ffffffe0, this is print No.1

switch to [PID = 3 COUNTER = 1 PRIORITY = 52]
[U-MODE] pid: 3, sp is 0000003ffffffe0, this is print No.1

switch to [PID = 2 COUNTER = 1 PRIORITY = 88]
[U-MODE] pid: 2, sp is 0000003ffffffe0, this is print No.1

switch to [PID = 1 COUNTER = 1 PRIORITY = 37]
[U-MODE] pid: 1, sp is 0000003ffffffe0, this is print No.1

switch to [PID = 2 COUNTER = 88 PRIORITY = 88]
[U-MODE] pid: 2, sp is 0000003ffffffe0, this is print No.2
[U-MODE] pid: 2, sp is 0000003ffffffe0, this is print No.3
QEMU: Terminated
fengnian@MICROSO-VU1H310: /mnt/d/study/OS/os23fall-stu/src/lab4$

```

与之前的结果相同

5. 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）

答：使用的是一对一的用户态线程和内核态线程

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

答：`p_filesz` 是段在可文件中空间大小，只包含了文件所占的实际字节数，`p_memsz` 是指段在内存中占的大小，它可能包含了未初始化的数据或者初始化为0的数据（比如.bss段，如果一个数据未被初始化，就不需要为其分配空间，所以.bss并不占用可执行文件的大小，仅仅记录需要用多少空间来存储这些未初始化的数据，而不分配实际空间），这些数据存储在磁盘上会浪费内存，只有在 ELF 文件加载到内存中后才会占用空间，因此占用了额外的内存，这就导致 `p_filesz` 的大小一般小于等于 `p_memsz`

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？

答：

1. 因为每个进程初始化的时候都申请了一块物理地址做栈，然后把虚拟地址映射过去，虽然多个进程的栈虚拟地址是相同的，但实际映射到物理地址上是不同的
2. 应该没有办法，这个物理地址只有内核态能看到