

浙江大学

本科实验报告

课程名称：操作系统

姓 名：徐文皓

学 院：计算机科学与技术学院

系：软件工程系

专 业：软件工程

学 号：3210102377

指导教师：夏莹杰

2023 年 11 月 20 日

浙江大学操作系统实验报告

实验名称：_____RV64 虚拟内存管理_____

电子邮件地址：_____手机：_____

实验地点：_____线上_____实验日期：2023 年 11 月 20 日

一、实验目的和要求

本实验的目的有：学习虚拟内存的相关知识，实现物理地址到虚拟地址的切换；了解 RISC-V 架构中 SV39 分页模式，实现虚拟地址到物理地址的映射，并对不同的段进行相应的权限设置。

本实验的要求是：独立完成作业；在 lab2 的基础上，修改 defs.h 文件；并且从仓库拉取 vmlinux.lds 和 Makefile 文件；进行虚拟内存映射，实现 setup_vm 和 setup_vm_final 函数。

二、实验过程

在 defs.h 中追加指定内容。

```
1 #define OPENSBI_SIZE (0x200000)
2
3 #define VM_START (0xffffffff00000000)
4 #define VM_END (0xfffffffff0000000)
5 #define VM_SIZE (VM_END - VM_START)
6
7 #define PA2VA_OFFSET (VM_START - PHY_START)
```

将本实验新增的文件添加到相应位置。在顶层的 Makefile 中做出如下更改，使得内核可以使用刷新缓存的指令扩展，并自动在编译项目前执行 clean 任务来防止对头文件的修改无法触发编译任务。

```
1 ISA=rv64imafd_zifencei
2 ...
3 all: clean
4 ...
```


由于 `vmlinux.lds` 被更改，在 `head.S` 中，经过 `gdb` 调试，我们发现此时 `sp` 已经被设置为虚拟地址。但由于此时 `setup_vm()` 等尚未被调用，因此我们先通过 `sp=sp-PA2VA_OFFSET` 将其设置回物理地址，对 `setup_vm()` 进行调用。返回后，物理地址到虚拟地址的第一次映射已完成，调用 `relocate`。

```

1 _start:
2   la sp, boot_stack_top
3   li t0, 0xfffffffff8000000 // PA2VA_OFFSET
4   sub sp, sp, t0
5
6   call setup_vm
7   call relocate
8   call mm_init

```

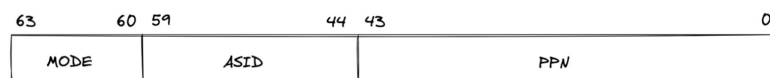
在 `relocate` 中，将 `ra` 和 `sp` 移动到虚拟地址，并对 `satp` 进行设置。

```

1 relocate:
2   # set ra = ra + PA2VA_OFFSET
3   # set sp = sp + PA2VA_OFFSET (If you have set the sp before)
4
5   li t0, 0xfffffffff8000000 // PA2VA_OFFSET
6   add ra, ra, t0
7   add sp, sp, t0
8
9   # set satp with early_pgtbl
10  la t0, early_pgtbl
11  li t1, 0xfffffffff8000000 // PA2VA_OFFSET
12  sub t0, t0, t1
13  srlj t0, t0, 12
14  li t1, 0x8000000000000000
15  or t0, t0, t1
16  csrw satp, t0
17
18  # flush tlb
19  sfence.vma zero, zero
20
21  # flush icache
22  fence.i
23
24  ret

```

考虑 `satp` 的结构。`satp[43:0]`(PPN)存放顶级页表(即 `early_pgtbl`)对应的物理页号，即 $(\text{early_pgtbl} - \text{PA2VA_OFFSET}) \gg 12$ 。`satp[59:44]`(ASID)在本次实验中置零。`satp[63:60]`(MODE)设置为 8，即选择 SV39 模式。



satp Register

从 `relocate` 返回后，`gdb` 显示跳转到高地址处继续执行。稍后在 `set_vm_final()` 中，我们完成对其余地址的映射。

在进行这一步之前，现在 `mm.c` 中修改 `mm_init` 函数，实现对内存的初始化。考虑 `kfreerange` 的函数原型为 `kfreerange(char *start, char *end)`，而我们观察到，在 `vmlinux.lds` 中，第一个实际参数 `_skernel` 的地址已经被修改到高地址，本次实验中由物理地址到虚拟地址的映射为线性映射，因此只需要将第二个参数修改为

(char *)PHY_END+PA2VA_OFFSET 即可。

```
1 /* kernel代码起始位置 */
2 BASE_ADDR = VM_START + OPENSBI_SIZE;
3 ...
4 /* . 代表当前地址 */
5 . = BASE_ADDR;
6 /* 记录kernel代码的起始地址 */
7 _skernel = .;
8 ...
```

```
1 void mm_init(void) {
2     kfree_range(_ekernel, (char *)PHY_END+PA2VA_OFFSET);
3     printk("...mm_init done!\n");
4 }
```

在 vmlinux.lds 中，可以观察到 _stext 和 _etext 分别代表 text 段的起始和结束地址，以此类推。而 _sdata 代表其余部分的起始地址。因此，我们将这些地址标识符声明到 vm.c 中，以便 setup_vm_final() 使用。

```
1 /* _stext, _etext 分别记录了text段的起始与结束地址 */
2 .text : ALIGN(0x1000){
3     _stext = .;
4     ...
5     _etext = .;
6 } >ramv AT>ram
7
8 .rodata : ALIGN(0x1000){
9     _srodata = .;
10    ...
11    _erodata = .;
12 } >ramv AT>ram
13
14 .data : ALIGN(0x1000){
15     _sdata = .;
16    ...
```

在 setup_vm_final 中，分别调用 create_mapping，对 text、rodata 和其余部分分别建立映射，起始地址和结束地址用刚刚引入的外部声明表示，长度为两者之差，并通过最后一个参数设置页权限。需要特别指明的是，对于其余部分，长度为 PHY_END+PA2VA_OFFSET-(uint64)_sdata。在完成映射后，将 satp 的 PPN 部分设置为 swapper_pg_dir 所对应的物理地址。

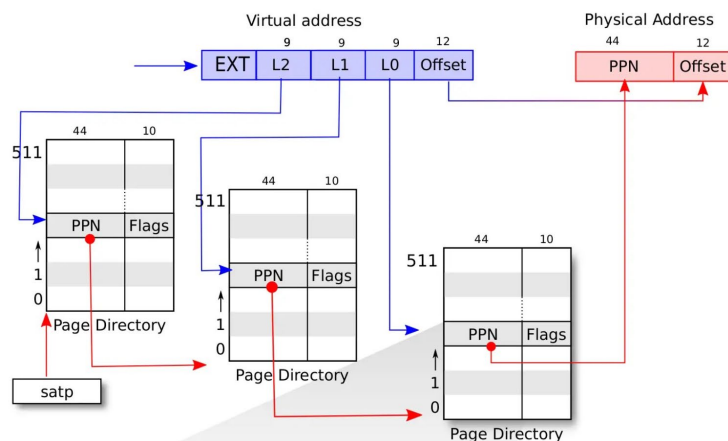
```
1 unsigned long swapper_pg_dir[512] __attribute__((__aligned__(0x1000)));
2 extern char _stext[], _etext[];
3 extern char _srodata[], _erodata[];
4 extern char _sdata[], _edata[];
5
6 void setup_vm_final(void) {
7     memset(swapper_pg_dir, 0x0, PGSIZE);
8
9     // mapping kernel text X[-|R|V
10    create_mapping(swapper_pg_dir, _stext, _stext-PA2VA_OFFSET, _etext-_stext, 0xb);
11
12    // mapping kernel rodata -[-|R|V
13    create_mapping(swapper_pg_dir, _srodata, _srodata-PA2VA_OFFSET, _erodata-_srodata, 0x3);
14
15    // mapping other memory -[W|R|V
16    create_mapping(swapper_pg_dir, _sdata, _sdata-PA2VA_OFFSET, PHY_END+PA2VA_OFFSET-(uint64)_sdata, 0x7);
17
18    // set satp with swapper_pg_dir
19
20    register unsigned long t0 asm("t0") = (unsigned long)swapper_pg_dir - PA2VA_OFFSET;
21    t0 = (t0 >> 12) | 0x8000000000000000;
22    asm volatile("csrw satp, t0" : : "r"(t0) : "memory");
23
24    // flush TLB
25    asm volatile("sfence.vma zero, zero");
26
27    // flush icache
28    asm volatile("fence.i");
29    return;
30 }
```

接下来考虑 `create_mapping()`。为便于操作，我们引入 `create_mapping_page()`，并对 `create_mapping()` 按页进行映射。由于 `va` 和 `pa` 均为地址而非页码，因此步长需要设置为 `PGSIZE`。

```
1 create_mapping(uint64 *pgtbl, uint64 va, uint64 pa, uint64 sz, int perm) {
2     for (int i = 0; i < sz; i = i + PGSIZE)
3         create_mapping_page(pgtbl, va + i, pa + i, perm);
4 }
```

在 `create_mapping_page()` 中，分别使用掩码取 `va` 的不同位，获得 3 个 VPN。首先查看 `pgtbl[VPN[2]]`，它对应第一级页表的 `pte`，通过 `0x1` 掩码得到 `valid` 位，如果为 0，说明该页表项还不存在，我们通过 `kalloc()` 获得一页作为页表目录，以其物理地址和 `0x1` 的 `XWRV` 设置 `Pagetable Entry`，作为 `pgtbl[VPN[2]]`。随后，通过 `pgtbl = (unsigned long *)((pgtbl[VPN[2]] >> 10 << 12) + PA2VA_OFFSET)` 来到这个页表，对下一级页表进行类似操作。在第三级页表时，它将指向物理地址而非下一个页表，因此我们将这个 `Pagetable Entry` 的 `XWRV` 设置为指定参数，当 `XWR` 不为 0 时说明该 `Pagetable Entry` 所指向的是物理地址。

```
1 create_mapping_page(uint64 *pgtbl, uint64 va, uint64 pa, int perm){
2     unsigned long VPN[3];
3     VPN[2] = (va >> 30) & 0x1fff;
4     VPN[1] = (va >> 21) & 0x1fff;
5     VPN[0] = (va >> 12) & 0x1fff;
6
7     unsigned long pte = pgtbl[VPN[2]]; // 第一级映射
8     if ((pte & 0x1) == 0) { // 如果不存在页表
9         unsigned long pgtbl_tmp = (unsigned long)kalloc() - PA2VA_OFFSET;
10        pgtbl_tmp = (pgtbl_tmp >> 12) << 10;
11        pgtbl[VPN[2]] = (unsigned long)pgtbl_tmp | 0x1;
12    }
13
14    pgtbl = (unsigned long *)((pgtbl[VPN[2]] >> 10 << 12) + PA2VA_OFFSET);
15
16    pte = pgtbl[VPN[1]]; // 第二级映射
17    if ((pte & 0x1) == 0) { // 如果不存在页表
18        unsigned long pgtbl_tmp = (unsigned long)kalloc() - PA2VA_OFFSET;
19        pgtbl_tmp = (pgtbl_tmp >> 12) << 10;
20        pgtbl[VPN[1]] = (unsigned long)pgtbl_tmp | 0x1;
21    }
22
23    pgtbl = (unsigned long *)((pgtbl[VPN[1]] >> 10 << 12) + PA2VA_OFFSET);
24
25    pte = ((pa >> 12) << 10) | perm;
26    pgtbl[VPN[0]] = pte;
27 }
```



至此,我们完成了 `setup_vm_final()` 的相关设置,在 `head.S` 中将其置于 `mm_init` 和 `task_init` 之间即可。

```
1 ...
2 call mm_init
3 call setup_vm_final
4 call task_init
5 ...
```

至此,本实验的全部任务已经完成。

为便于观察实验现象,我们修改 `arch/riscv/kernel/proc.c`,在 `dummy()` 函数中将原来的输出信息追加上进程地址。

```
1 printk("[PID = %d] is running. thread space begins at 0x%llx\n", current->pid, current);
```

启动内核,可以观察到以下现象。

```
...proc_init done!
2022 Hello RISC-V
[S] Supervisor Mode Timer Interrupt

SWITCH TO [pid=8 counter=1 priority=25]
[PID = 8] is running. thread space begins at 0xffffffff007fb7000
[S] Supervisor Mode Timer Interrupt

SWITCH TO [pid=7 counter=2 priority=5]
[PID = 7] is running. thread space begins at 0xffffffff007fb8000
[S] Supervisor Mode Timer Interrupt
[PID = 7] is running. thread space begins at 0xffffffff007fb8000
[S] Supervisor Mode Timer Interrupt

SWITCH TO [pid=14 counter=2 priority=6]
[PID = 14] is running. thread space begins at 0xffffffff007fb1000
[S] Supervisor Mode Timer Interrupt
[PID = 14] is running. thread space begins at 0xffffffff007fb1000
[S] Supervisor Mode Timer Interrupt
```

三、讨论和心得

本次实验理解起来较有难度,尤其是对内存的操作使得 `gdb` 并不能像前几次实验那样正常工作,但是在同学的热心帮助下实验得以较为顺利地顺利完成。实验中,可以通过 `gdb` 输出相应变量、函数的地址,检查其位置是否处于虚拟地址。过程中遇到了以下问题:

在外部声明 `_stext` 等值时,我起初使用 `char*` 类型,无法正常运行,经过 `gdb` 调试,发现 `&_stext` 才是其地址。将声明改为 `char[]` 类型或者将 `create_mapping()` 相应实际参数更改为形如 `&_stext` 的格式,问题得以解决。这里启示我们在对值进行强制类型转换时,转换到指针类型和转换到数组类型的区别。

建议在实验指导书中对 `vmlinux.lds` 的更改做一些说明。

四、思考题

1. 验证.text, .rodata 段的属性是否成功设置，给出截图。

即验证.text 段的 XWR 是否为 101，验证.rodata 段的 XWR 是否为 001。

程序可以正常执行，说明.text 段的 X 位成功设置为 1。

要验证某个段是否具有某种权限，只需要验证在对该段进行某种操作时是否会抛出异常即可。在之前的实验中，我们已经在中断处理函数中对非时钟中断输出 trap 原因，可以借此来进行观察。

```
1 //arch/riscv/kernel/trap.c
2 ...
3 if(interrupt&&cause==5){ // timer interrupt
4     printk("[S] Supervisor Mode Timer Interrupt\n");
5     clock_set_next_event();
6     do_timer();
7 }
8 else{
9     printk("[Error] Not Timer Interrupt with scause = 0x%16x\n", scause);
10 }
11 ...
```

然而，非时钟中断会导致程序一直返回同一行进而出现死循环。为了打破这个局面，我们在处理其他中断时对 sepc 加 4。考虑到 entry.S 中是由 _traps 在调用 trap_handler() 前后负责维护 sepc 的稳定性，因此在 trap_handler() 中实现这一操作并不可行。

所以，我们在 _traps 中实现这一功能：在从 trap_handler 返回后，如果 scause 表明本次中断是时钟中断则将 sepc 恢复，否则将 sepc 恢复后加 4。

```
1 # arch/riscv/kernel/entry.S
2 ...
3     call trap_handler
4
5     # 3. restore sepc and 32 registers (x2(sp) should be restore last) from stack
6     ld a0, 256(sp)
7     li t0, 0x8000000000000005 # temporarily add 4 to sepc
8     csrr t1, scause
9     beq t1, t0, _csrwrite
10    addi a0, a0, 4
11 _csrwrite:
12     csrw sepc, a0
13 ...
```

至此，我们可以通过观察异常信息检查两个段的权限。

首先检验.rodata 段的 X 权限。在 head.S 中，我们完成中断设置后，执行 j _srodata 指令，即试图跳转到.rodata 段进行执行。

```
1 # arch/riscv/kernel/head.S
2     csrs sstatus, 2 # set sstatus[SIE] = 1
3
4     j _srodata
5     call start_kernel
```


启动内核，发现给出了 `scause` 为 `0xc` 的异常提示。

```
...proc_init done!
[Error] Not Timer Interrupt with scause = 0x0000000c
```

对照 `scause` 的表，我们发现 `0xc` 为 `Instruction Page Fault`，说明内核无法从 `__srodata` 处读取指令，即不能执行 `.rodata` 的内容。

Interrupt	Exception Code	Description
1	0	User software interrupt
1	1	Supervisor software interrupt
1	2-3	Reserved for future standard use
1	4	User timer interrupt
1	5	Supervisor timer interrupt
1	6-7	Reserved for future standard use
1	8	User external interrupt
1	9	Supervisor external interrupt
1	10-15	Reserved for future standard use
1	≥16	Reserved for platform use
0	0	Instruction address misaligned
0	1	Instruction access fault
0	2	Illegal instruction
0	3	Breakpoint
0	4	Load address misaligned
0	5	Load access fault
0	6	Store/AMO address misaligned
0	7	Store/AMO access fault
0	8	Environment call from U-mode
0	9	Environment call from S-mode
0	10-11	Reserved for future standard use
0	12	Instruction page fault
0	13	Load page fault
0	14	Reserved for future standard use
0	15	Store/AMO page fault
0	16-23	Reserved for future standard use
0	24-31	Reserved for custom use
0	32-47	Reserved for future standard use
0	48-63	Reserved for custom use
0	≥64	Reserved for future standard use

接下来验证 `W` 和 `R` 两个属性。

在 `start_kernel()` 中，我们尝试打印 `__stext` 的值，并尝试对其写入数据。

```
1 // init/main.c
2 ...
3 extern char __stext[];
4
5 int start_kernel() {
6     ...
7
8     printk("__stext = %x\n", *__stext);
9     *__stext = 0;
10    printk("__stext = %x\n", *__stext);
11    ...
12
13 }
```

启动内核，观察到如下输出。能够正常打印其值，证明它是可读的。`scause` 为 `0xf`，对应 `Store/AMO Page Fault`，说明写操作被拒绝。相应地，在下次打印中我们发现其值并没有发生变化，这同样说明了写操作失败。

```
...proc_init done!
2022 Hello RISC-V
__stext = 00000017
[Error] Not Timer Interrupt with scause = 0x0000000f
__stext = 00000017
[S] Supervisor Mode Timer Interrupt
```

类似地，对 `__srodata` 进行验证，现象相同。

```
...proc_init done!
2022 Hello RISC-V
__srodata = 0000002e
[Error] Not Timer Interrupt with scause = 0x0000000f
__srodata = 0000002e
[S] Supervisor Mode Timer Interrupt
```

至此，我们验证了 `.text` 和 `.rodata` 段 `XWR` 属性的正确性。

2. 为什么我们在 `setup_vm()` 中需要做等值映射?

在回答这个问题前，不妨将等值映射删去，观察现象。

```
1 // arch/riscv/kernel/vm.c
2
3 // early_pgtbl[(va & mid_bits) >> 30] = ((pa >> 12) << 10) | 0xf;
```

启动内核，意料之中地，无法正常运行。经过 `gdb` 调试，我们发现在 `csrw satp,t0` 时程序无法继续执行。这里的功能是将 `satp` 的 PPN 位设置为 `early_pgtbl` 的物理地址。

```
0x8020009c      addiw   t1,zero,-1
0x802000a0      slliw   t1,t1,0x3f
> 0x802000a4      or      t0,t0,t1
0x802000a8      csw     satp,t0
0x802000ac      sfence.vma
0x802000b0      fence.i
0x802000b4      ret
```

在调用 `setup_vm_final()` 完成对所有页表项的设置之前，我们的内核还在 `*0x8020000` 处运行。在修改 `satp` 后、`ret` 执行前程序仍在物理地址处运行，如果在此时不进行等值映射，随后的命令(在本例中，对应 `*0x802000ac` 的 `sfence.vma` 指令，它已经认为这个地址是虚拟地址而非物理地址)将无法从新的 `satp` 所指向的页表中查找到实际的物理地址，进而无法执行。

为了验证这里是由于命令无法执行所导致程序无法正常运行，我们在等值映射时通过修改参数改变该页的权限。当参数为 `0xf` 时，即实验中的情况，我们已经知道可以正常运行。将参数修改为 `0x7`，即授予 `WR` 权限而不授予 `X` 权限，仍未在 `csrw satp,t0` 命令处无法继续执行，问题并没有解决。将参数修改为 `0x9`，即授予 `X` 权限而不授予 `WR` 权限，可以发现程序可以正常运行。

因此，我们判断是 `satp` 发生变更后的三条指令仍需要访问页表确定自己的物理地址，进而被执行。而等值映射为这三条指令能够“找到”自己的物理地址提供了保证。

3. 在 Linux 中，是不需要做等值映射的。请探索一下不在 `setup_vm()` 中做等值映射的方法。

经由上文中我们所发现的程序正常运行与该页 X 权限的依赖关系这一事实，我们可以发现在 `satp` 发生变更后，程序在执行下一条指令 `sfence.vma` 时，会触发 `Instruction Page Fault` 这一异常。考虑 `trap` 处理函数，这将是我们的突破口。

然而，如果能够正常进入 `trap` 处理函数，根据我们 `trap_handler()` 的设计逻辑，我们所观察到的现象应该是时钟中断和其他 `trap` 中的某一种输出。程序呈现卡死状态，说明这一函数并没有顺利执行。经过思考，我们发现问题在于 `stvec` 受到了虚拟地址的影响，无法从新的 `satp` 所指向的页表中查找到实际的物理地址。

为此，不妨考虑更新 `stvec` 的值，使其能够被正常访问。在 `head.S` 的 `relocate` 中，设置 `satp` 前，我们临时修改 `stvec` 的值，使得 `satp` 设置后程序产生异常能够跳转到 `handle_setup_vm`，在这里我们讲 `sepc` 调整至正确地址。

```
1 relocate:
2 ...
3 la t1, handle_setup_vm
4 csrw stvec, t1
5
6 csrw satp, t0
7
8 # flush tlb
9 sfence.vma zero, zero
10
11 # flush icache
12 fence.i
13
14 ret
15
16 handle_setup_vm:
17 csrr t1, sepc
18 li t2, 0xfffffffdf8000000 // PA2VA_OFFSET
19 add t1, t1, t2
20 addi t1, t1, -4
21 csrw sepc, t1
22 sret
```

需要特别说明的是，异常与中断不同，发生异常时 `sepc` 的值会被设置为异常指令的下一条指令地址，但在本例中我们需要发生异常的 `sfence.vma` 语句完成执行，因此，在设置 `sepc` 时，我们引入了 -4 这一偏移，即向前一条指令。

下面给出实验验证：

在执行到临时修改 stvec 的位置，我们发现 sepc 的值为 0x802000bc。

```
head.S
68     ret
69
70 handle_setup_vm:
B+ 71     csrr t1, sepc
72     li t2, 0xfffffffff8000000 // PA2VA_OFFSET
73     add t1, t1, t2
> 74     addi t1, t1, -8
75     csrw sepc, t1
76     sret

B+ 0xffffffe0002000c4 <handle_setup_vm> csrr t1,sepc
0xffffffe0002000c8 <handle_setup_vm+4> addiw t2,zero,-65
0xffffffe0002000cc <handle_setup_vm+8> slli t2,t2,0x1f
0xffffffe0002000d0 <handle_setup_vm+12> add t1,t1,t2
> 0xffffffe0002000d4 <handle_setup_vm+16> addi t1,t1,-8
0xffffffe0002000d8 <handle_setup_vm+20> csrw sepc,t1
0xffffffe0002000dc <handle_setup_vm+24> sret
0xffffffe0002000e0 <_traps> addi sp,sp,-264
0xffffffe0002000e4 <_traps+4> sd zero,0(sp)
0xffffffe0002000e8 <_traps+8> sd ra,8(sp)

remote Thread 1.1 In: handle_setup_vm L74 PC: 0xffffffe0002000d4
$1 = 0xffffffe0002000c4
(gdb) b *0xffffffe0002000c4
Breakpoint 3 at 0xffffffe0002000c4: file head.S, line 71.
(gdb) c
Continuing.

Breakpoint 3, handle_setup_vm () at head.S:71
(gdb) n
(gdb) p/x $sepc
$2 = 0x802000bc
```

而我们预期在异常返回后的地址为 0x802000b8。

```
0x802000a0 slli t1,t1,0x3f
0x802000a4 or t0,t0,t1
0x802000a8 auipc t1,0x3
0x802000ac ld t1,-152(t1)
> 0x802000b0 csrw stvec,t1
0x802000b4 csrw satp,t0
0x802000b8 sfence.vma
0x802000bc fence.i
0x802000c0 ret
0x802000c4 csrr t1,sepc
```

因此，-4 偏移是合理的。

这样，内核即可在不作等值映射的条件下正常运行。

```
...proc_init done!
2022 Hello RISC-V
[S] Supervisor Mode Timer Interrupt

SWITCH TO [pid=8 counter=1 priority=25]
[PID = 8] is running. thread space begins at 0xffffffe007fb7000
[S] Supervisor Mode Timer Interrupt

SWITCH TO [pid=7 counter=2 priority=5]
[PID = 7] is running. thread space begins at 0xffffffe007fb8000
[S] Supervisor Mode Timer Interrupt
[PID = 7] is running. thread space begins at 0xffffffe007fb8000
[S] Supervisor Mode Timer Interrupt

SWITCH TO [pid=14 counter=2 priority=6]
[PID = 14] is running. thread space begins at 0xffffffe007fb1000
[S] Supervisor Mode Timer Interrupt
[PID = 14] is running. thread space begins at 0xffffffe007fb1000
```