

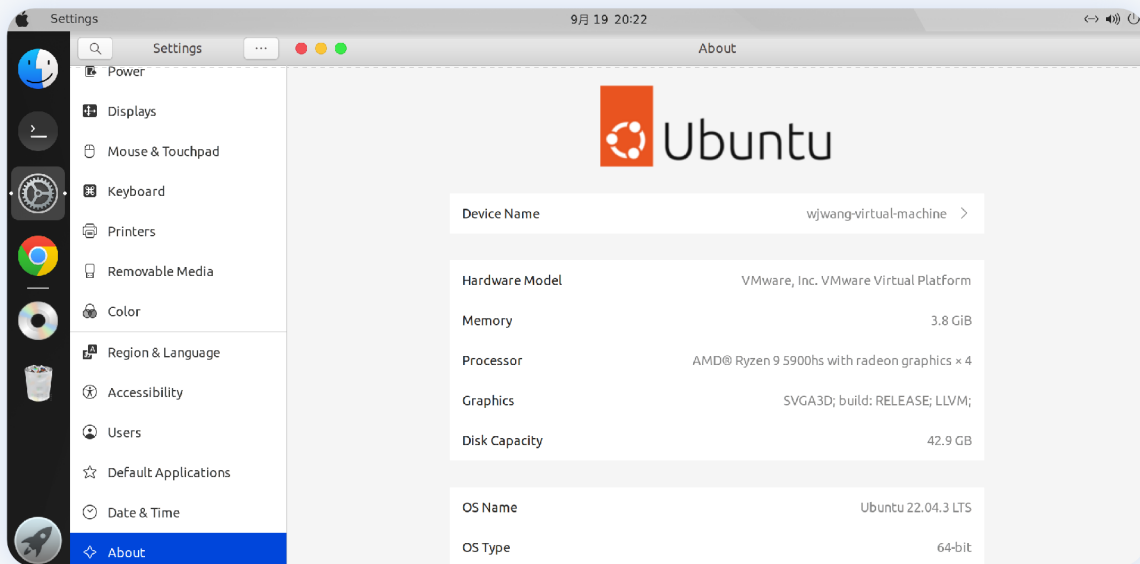
# Lab 3 RV64 虚拟内存管理

学号：3210106034

姓名：王伟杰

实验环境：

```
1  uname -a
2  Linux wjwang-virtual-machine 6.2.0-33-generic #33~22.04.1-Ubuntu SMP
    PREEMPT_DYNAMIC Thu Sep  7 10:33:52 UTC 2 x86_64 x86_64 x86_64
    GNU/Linux
```



## 准备工程

在 `defs.h` 添加如下内容

```

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

```

从 `repo` 同步 `vmlinux.lds`，并正确放置。

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  # Makefile
2  # ...
3  ISA=rv64imafd_zifencei
4  # ...
5  all: clean
6      ${MAKE} -C lib all
7      ${MAKE} -C test all
8      ${MAKE} -C init all
9      ${MAKE} -C arch/riscv all
10     @echo -e '\n'Build Finished OK
11  # ...

```

## 开启虚拟内存映射

### setup\_vm 的实现

对 `PHY_START` 开始的1GB区域进行两次映射：一次等值映射与一次映射到 `direct mapping area`。

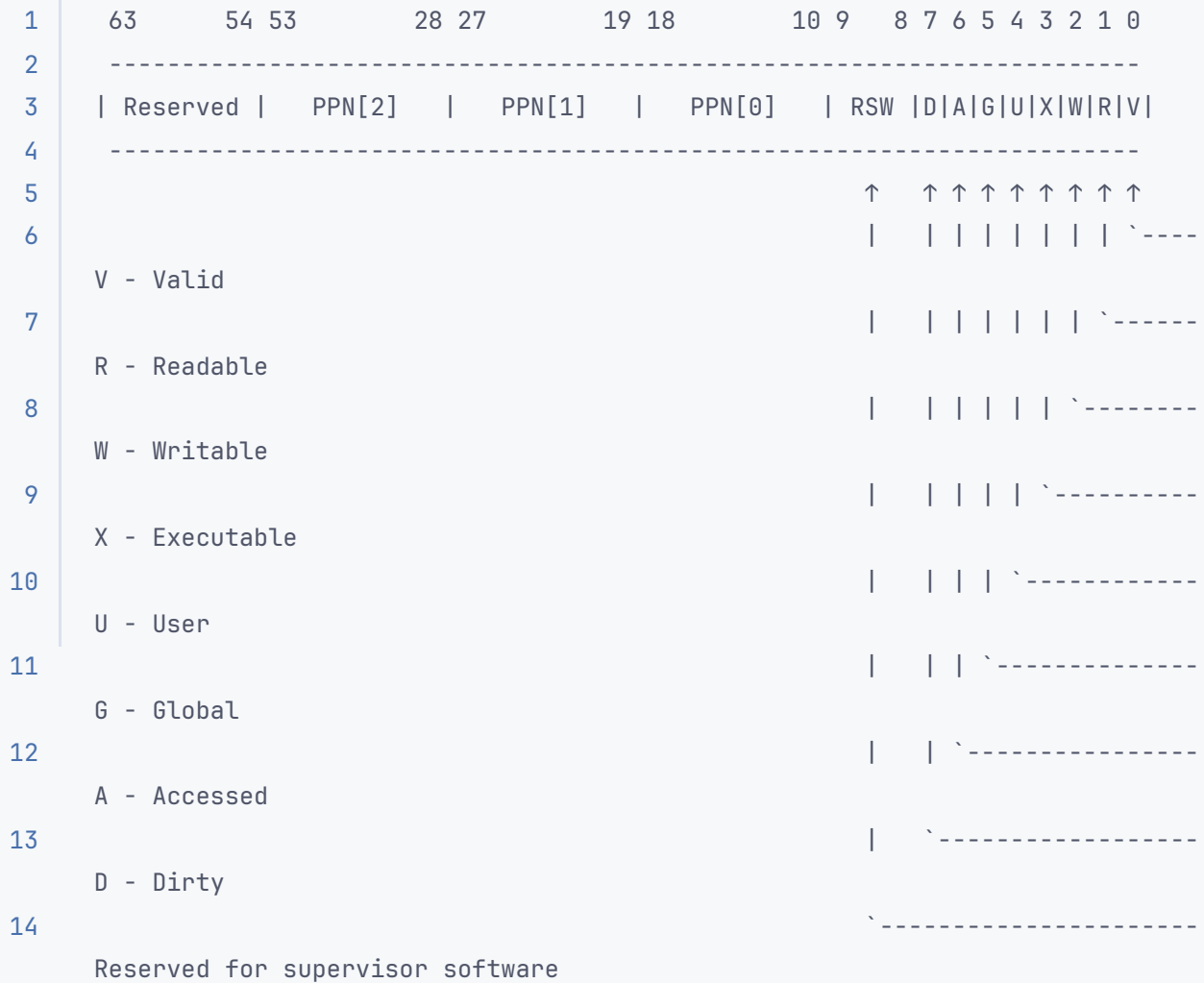
由于 `PA` 的格式如下：

```

1  55          30 29          21 20          12 11          0
2  -----
3  |      PPN[2]      |  PPN[1]  |  PPN[0]  |      page offset
4  |
5  -----
   Sv39 physical address

```

我们不需要使用多级页表，所以 PPN 的值即为上图中的 PPN[2]，PTE 格式如下：



所以我们将得到的 PPN 左移28位，再设置权限位即可。

```
1 // arch/riscv/kernel/vm.c
2 void setup_vm(void) {
3     /*
4         1. 由于是进行 1GB 的映射 这里不需要使用多级页表
5         2. 将 va 的 64bit 作为如下划分： | high bit | 9 bit | 30 bit |
6             high bit 可以忽略
7             中间9 bit 作为 early_pgtbl 的 index
8             低 30 bit 作为 页内偏移 这里注意到 30 = 9 + 9 + 12， 即我们只使用根页表， 根页
9             表的每个 entry 都对应 1GB 的区域。
10        3. Page Table Entry 的权限 V | R | W | X 位设置为 1
11    */
12    unsigned long PA = PHY_START;
13    unsigned long VA_EQ = PA;
14    int index = (VA_EQ >> 30) & 0x1fff;
```

```

14     unsigned long PPN = PA >> 30 & 0x3fffffff;
15     unsigned long PTE = (PPN << 28) | 0xf; // V R W X = 1
16     early_pgtbl[index] = PTE;
17
18     unsigned long VA_DIRECT = PA + PA2VA_OFFSET;
19     index = (VA_DIRECT >> 30) & 0x1ff;
20     PPN = PA >> 30 & 0x3fffffff;
21     PTE = (PPN << 28) | 0xf;
22     early_pgtbl[index] = PTE;
23
24     printk("...setup_vm done!\n");
25 }

```

这里还需要更改设置 `sp` 到 `boot_stack` 的顶部的代码，因为 `vmlinux.lds` 将kernel代码起始位置修改为了 `VM_START + OPENSBI_SIZE`：

```

1  # arch/riscv/kernel/head.S
2  # ...
3  _start:
4      # -----
5      # - your code here -
6      # 设置sp到boot_stack的顶部
7      li t0, 0xffffffffdf80000000
8      la sp, boot_stack_top
9      sub sp, sp, t0
10
11     call setup_vm
12     call relocate
13
14     call mm_init
15     call task_init
16     # ...

```

`satp` 寄存器 ( `PA >> 12 = PPN` )：

1	63	60 59	44 43	0
2	-----			
3		MODE		ASID
4	-----			PPN
	-----			

同时我们完成 `relocate` 函数，完成对 `satp` 的设置，以及跳转到对应的虚拟地址：

```
1  relocate:
2      # set ra = ra + PA2VA_OFFSET
3      # set sp = sp + PA2VA_OFFSET (If you have set the sp before)
4      # PA2VA_OFFSET = 0xffffffffe000000000 - 0x0000000080000000 =
0xffffffffdf80000000
5      li t1, 0xffffffffdf80000000
6      add ra, ra, t1
7      add sp, sp, t1
8
9      # set satp with early_pgtbl
10     la t0, early_pgtbl
11     sub t0, t0, t1
12     srl t0, t0, 12
13     li t1, 0x8000000000000000
14     or t0, t0, t1
15     csrw satp, t0
16
17     # flush tlb
18     sfence.vma zero, zero
19
20     # flush icache
21     fence.i
22
23     ret
```

## setup\_vm\_final 的实现

修改 `mm.c` 中的代码，`mm.c` 中初始化的函数接收的起始结束地址调整为虚拟地址：

```
1  void mm_init(void) {
2      kfreerange(_kernel, (char *)VM_START + PHY_SIZE);
3      printk("...mm_init done!\n");
4  }
```

由于之前映射的页表比较粗糙，并且后续不再需要等值映射，所以我们在这里对所有物理内存 (128M) 进行映射，并设置正确的权限：

```
1  // arch/riscv/kernel/vm.c
```

```

2  /* swapper_pg_dir: kernel pagetable 根目录, 在 setup_vm_final 进行映射。 */
3  unsigned long swapper_pg_dir[512] __attribute__((__aligned__(0x1000)));
4  extern char _stext[], _etext[];
5  extern char _srodata[], _erodata[];
6  extern char _sdata[], _edata[];
7  void create_mapping(uint64 *pgtbl, uint64 va, uint64 pa, uint64 sz, uint64
perm);

8
9  void setup_vm_final(void) {
10     memset(swapper_pg_dir, 0x0, PGSIZE);
11
12     // No OpenSBI mapping required
13
14     // mapping kernel text X|→R|V
15     create_mapping((uint64 *)swapper_pg_dir, (uint64)_stext, (uint64)_stext -
PA2VA_OFFSET, (uint64)(_etext - _stext), 0xb);
16
17     // mapping kernel rodata →→R|V
18     create_mapping((uint64 *)swapper_pg_dir, (uint64)_srodata,
(uint64)_srodata - PA2VA_OFFSET, (uint64)(_erodata - _srodata), 0x3);
19
20     // mapping other memory →W|R|V
21     create_mapping((uint64 *)swapper_pg_dir, (uint64)_sdata, (uint64)_sdata -
PA2VA_OFFSET, PHY_END + PA2VA_OFFSET - (uint64)_sdata, 0x7);
22
23     // set satp with swapper_pg_dir
24     unsigned long satp = 0x8000000000000000 | ((unsigned long)swapper_pg_dir
- PA2VA_OFFSET) >> 12;
25     csr_write(satp, satp);
26
27     // flush TLB
28     asm volatile("sfence.vma zero, zero");
29
30     // flush icache
31     asm volatile("fence.i");
32     return;
33 }

```

在 `create_mapping` 中设置三级页表映射关系:

```

1  // arch/riscv/kernel/vm.c

```

```

2 void create_mapping(uint64 *pgtbl, uint64 va, uint64 pa, uint64 sz, uint64
perm) {
3     /*
4     pgtbl 为根页表的基地址
5     va, pa 为需要映射的虚拟地址、物理地址
6     sz 为映射的大小, 单位为字节
7     perm 为映射的权限 (即页表项的低 8 位)
8
9     创建多级页表的时候可以使用 kalloc() 来获取一页作为页表目录
10    可以使用 V bit 来判断页表项是否存在
11    */
12    int page_num = (sz + PGSIZE - 1) / PGSIZE; // 取整
13    for (int i = 0; i < page_num; i++) {
14        uint64 VPN[3];
15        VPN[2] = (va >> 30) & 0x1ff; // 9 bit
16        VPN[1] = (va >> 21) & 0x1ff;
17        VPN[0] = (va >> 12) & 0x1ff;
18        uint64 *pte = pgtbl;
19        for (int j = 2; j > 0; j--) {
20            if ((pte[VPN[j]] & 0x1) == 0) { // 如果valid为0, 则需要开辟新的pte
21                uint64 new_pte = kalloc();
22                pte[VPN[j]] = (((new_pte - PA2VA_OFFSET) >> 12) &
0xffffffffffff) << 10 | 0x1;
23                pte = (uint64 *)new_pte;
24            }
25            else {
26                pte = (uint64 *)((pte[VPN[j]] >> 10 << 12) + PA2VA_OFFSET);
27            }
28        }
29        pte[VPN[0]] = (((pa >> 12) & 0xffffffffffff) << 10) | perm;
30        va += PGSIZE;
31        pa += PGSIZE;
32    }
33 }

```

## 编译及测试

```
Boot HART Features      : scounteren,mcounteren,time
Boot HART PMP Count     : 16
Boot HART PMP Granularity : 4
Boot HART PMP Address Bits: 54
Boot HART MHPM Count    : 0
Boot HART MHPM Count    : 0
Boot HART MIDELEG       : 0x00000000000000222
Boot HART MEDELEG       : 0x0000000000000b109
...setup_vm done!
...mm_init done!
...proc_init done!
2022 Hello RISC-V
sstatus = 80000000000006002
sscratch = 0
[INTERRUPT] S mode timer interrupt!

switch to [PID = 8 COUNTER = 1 PRIORITY = 25]
[PID = 8] is running. auto_inc_local_var = 1
[INTERRUPT] S mode timer interrupt!

switch to [PID = 17 COUNTER = 1 PRIORITY = 71]
[PID = 17] is running. auto_inc_local_var = 1
[INTERRUPT] S mode timer interrupt!

switch to [PID = 18 COUNTER = 1 PRIORITY = 64]
[PID = 18] is running. auto_inc_local_var = 1
[INTERRUPT] S mode timer interrupt!

switch to [PID = 7 COUNTER = 2 PRIORITY = 5]
[PID = 7] is running. auto_inc_local_var = 1
[INTERRUPT] S mode timer interrupt!
[PID = 7] is running. auto_inc_local_var = 2
[INTERRUPT] S mode timer interrupt!

switch to [PID = 14 COUNTER = 2 PRIORITY = 6]
[PID = 14] is running. auto_inc_local_var = 1
[INTERRUPT] S mode timer interrupt!
[PID = 14] is running. auto_inc_local_var = 2
[INTERRUPT] S mode timer interrupt!
```

## 思考题

1. 验证 `.text` , `.rodata` 段的属性是否成功设置, 给出截图。

在成功设置之后, 我们知道 `.text` , `.rodata` 段是不允许写入的, 所以在设置之后对其进行写入, 如果可以成功写入则证明没有设置成功。

在 `setup_vm_final` 函数的末尾加入如下代码, 尝试对 `.text` , `.rodata` 段写入:



```

64
65     printk("...setup_vm_final done!\n");
66
67     printk("_stext = %x\n", *_stext);
68     printk("_srodata = %x\n", *_srodata);
69
70     *_stext = 0x00;
71     printk("*_stext = %x\n", *_stext);
72     *_srodata = 0x00;
73     printk("*_srodata = %x\n", *_srodata);
74
75     return;
76

```

先注释掉对 `.rodata` 进行操作的部分，运行代码，由于权限不够所以程序暂停运行。

```

make run
Domain0 HARTs          : 0*
Domain0 Region00       : 0x00000000080000000-0x0000000008001ffff ( )
Domain0 Region01       : 0x00000000000000000-0xfffffffffffff (R,W,X)
Domain0 Next Address   : 0x00000000080200000
Domain0 Next Arg1      : 0x00000000087000000
Domain0 Next Mode      : S-mode
Domain0 SysReset       : yes

Boot HART ID           : 0
Boot HART Domain       : root
Boot HART ISA          : rv64imafdcsu
Boot HART Features     : scouteren,mcounteren,time
Boot HART PMP Count    : 16
Boot HART PMP Granularity : 4
Boot HART PMP Address Bits: 54
Boot HART MHPM Count   : 0
Boot HART MHPM Count   : 0
Boot HART MIDELEG      : 0x00000000000000222
Boot HART MEDELEG      : 0x0000000000000b109
...setup_vm done!
...mm_init done!
...setup_vm_final done!
_stext = 0000009b

```

同理，注释掉对 `.text` 段的操作：

```
make run
Domain0 HARTs      : 0*
Domain0 Region00   : 0x0000000080000000-0x000000008001ffff ()
Domain0 Region01   : 0x0000000000000000-0xffffffffffff (R,W,X)
Domain0 Next Address : 0x0000000080200000
Domain0 Next Arg1   : 0x0000000087000000
Domain0 Next Mode   : S-mode
Domain0 SysReset    : yes

Boot HART ID       : 0
Boot HART Domain   : root
Boot HART ISA      : rv64imafdcsv
Boot HART Features : scounteren,mcounteren,time
Boot HART PMP Count : 16
Boot HART PMP Granularity : 4
Boot HART PMP Address Bits: 54
Boot HART MHPM Count : 0
Boot HART MHPM Count : 0
Boot HART MIDELEG  : 0x0000000000000222
Boot HART MEDELEG  : 0x000000000000b109
...setup_vm done!
...mm_init done!
...setup_vm_final done!
_srodata = 0000002e
```

证明 `.text` , `.rodata` 段的属性已经成功设置。

## 2.为什么我们在 `setup_vm` 中需要做等值映射?

因为在 `relocate` 执行过程中, 虽然会设置 `satp` 寄存器的值, 但是在设置完之后, `pc` 仍然是物理地址, 所以需要在原本的位置做一次等值映射, 保证这个函数仍然可以运行这个位置的代码, 使之顺利运行完。由于我们设置过 `ra` , 所以这个函数运行完之后会返回到虚拟地址继续运行。

## 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。

查找Linux源码, Linux是这样写`relocate`函数的:

```
1 relocate:
2     /* Relocate return address */
3     li a1, PAGE_OFFSET
4     la a0, _start
5     sub a1, a1, a0
6     add ra, ra, a1
7
8     /* Point stvec to virtual address of instruction after satp write */
9     la a0, 1f
10    add a0, a0, a1
11    csrw stvec, a0
12
```

```

13      /* Compute satp for kernel page tables, but don't load it yet */
14      la a2, swapper_pg_dir
15      srl a2, a2, PAGE_SHIFT
16      li a1, SATP_MODE
17      or a2, a2, a1
18
19      /*
20       * Load trampoline page directory, which will cause us to trap to
21       * stvec if VA  $\neq$  PA, or simply fall through if VA = PA
22       */
23      la a0, trampoline_pg_dir
24      srl a0, a0, PAGE_SHIFT
25      or a0, a0, a1
26      sfence.vma
27      csrw sptbr, a0
28  1:
29      /* Set trap vector to spin forever to help debug */
30      la a0, .Lsecondary_park
31      csrw stvec, a0

```

从中我们获得灵感，可以从 `exception` 入手调整 `pc` 在向虚拟地址转换过程中找不到正确路径的问题，首先我们删除 `setup_vm` 中等值映射的代码：

```

1  void setup_vm(void) {
2      /*
3       1. 由于是进行 1GB 的映射 这里不需要使用多级页表
4       2. 将 va 的 64bit 作为如下划分： | high bit | 9 bit | 30 bit |
5          high bit 可以忽略
6          中间9 bit 作为 early_pgtbl 的 index
7          低 30 bit 作为 页内偏移 这里注意到 30 = 9 + 9 + 12， 即我们只使用根页表，
          根页表的每个 entry 都对应 1GB 的区域。
8       3. Page Table Entry 的权限 V | R | W | X 位设置为 1
9      */
10     unsigned long PA, VA_EQ, VA_DIRECT, PPN, PTE;
11     int index;
12     PA = PHY_START;
13     // VA_EQ = PA;
14     // index = (VA_EQ >> 30) & 0x1fff;
15     // PPN = PA >> 30 & 0x3fffffff;
16     // PTE = (PPN << 28) | 0xf; // V R W X = 1
17     // early_pgtbl[index] = PTE;
18

```

```

19     VA_DIRECT = PA + PA2VA_OFFSET;
20     index = (VA_DIRECT >> 30) & 0x1ff;
21     PPN = PA >> 30 & 0x3ffffff;
22     PTE = (PPN << 28) | 0xf;
23     early_pgtbl[index] = PTE;
24
25     printk("...setup_vm done!\n");
26 }

```

然后在 `relocate` 函数中加入处理 `pc` 的错误跳转，在转换时trap to `stvec` 到标签为 `1` 的位置：

```

1  relocate:
2      # relocate return address
3      li t1, 0xffffffffdf80000000
4      add ra, ra, t1
5      add sp, sp, t1
6
7      # point stvec to virtual address of instruction after satp write
8      la t0, 1f
9      csrw stvec, t0
10
11     # set satp with early_pgtbl
12     la t0, early_pgtbl
13     sub t0, t0, t1
14     srl t0, t0, 12
15     li t1, 0x8000000000000000
16     or t0, t0, t1
17     csrw satp, t0
18
19
20     # flush tlb
21     sfence.vma zero, zero
22
23     # flush icache
24     fence.i
25
26     ret
27
28 1:
29     li t0, 0xffffffffdf80000000
30     csrr t1, sepc

```

```
31      add t1, t1, t0
32      csrw sepc, t1
33      ret
```

运行结果如下，可以正常编译运行：

```
Boot HART MHPM Count      : 0
Boot HART MIDELEG         : 0x00000000000000222
Boot HART MEDELEG         : 0x0000000000000b109
...setup_vm done!
...mm_init done!
...setup_vm_final done!
...proc_init done!
2022 Hello RISC-V
sstatus = 80000000000006102
sscratch = 0
[INTERRUPT] S mode timer interrupt!

switch to [PID = 8 COUNTER = 1 PRIORITY = 25]
[PID = 8] is running. auto_inc_local_var = 1
[INTERRUPT] S mode timer interrupt!

switch to [PID = 17 COUNTER = 1 PRIORITY = 71]
[PID = 17] is running. auto_inc_local_var = 1
[INTERRUPT] S mode timer interrupt!

switch to [PID = 18 COUNTER = 1 PRIORITY = 64]
[PID = 18] is running. auto_inc_local_var = 1
```