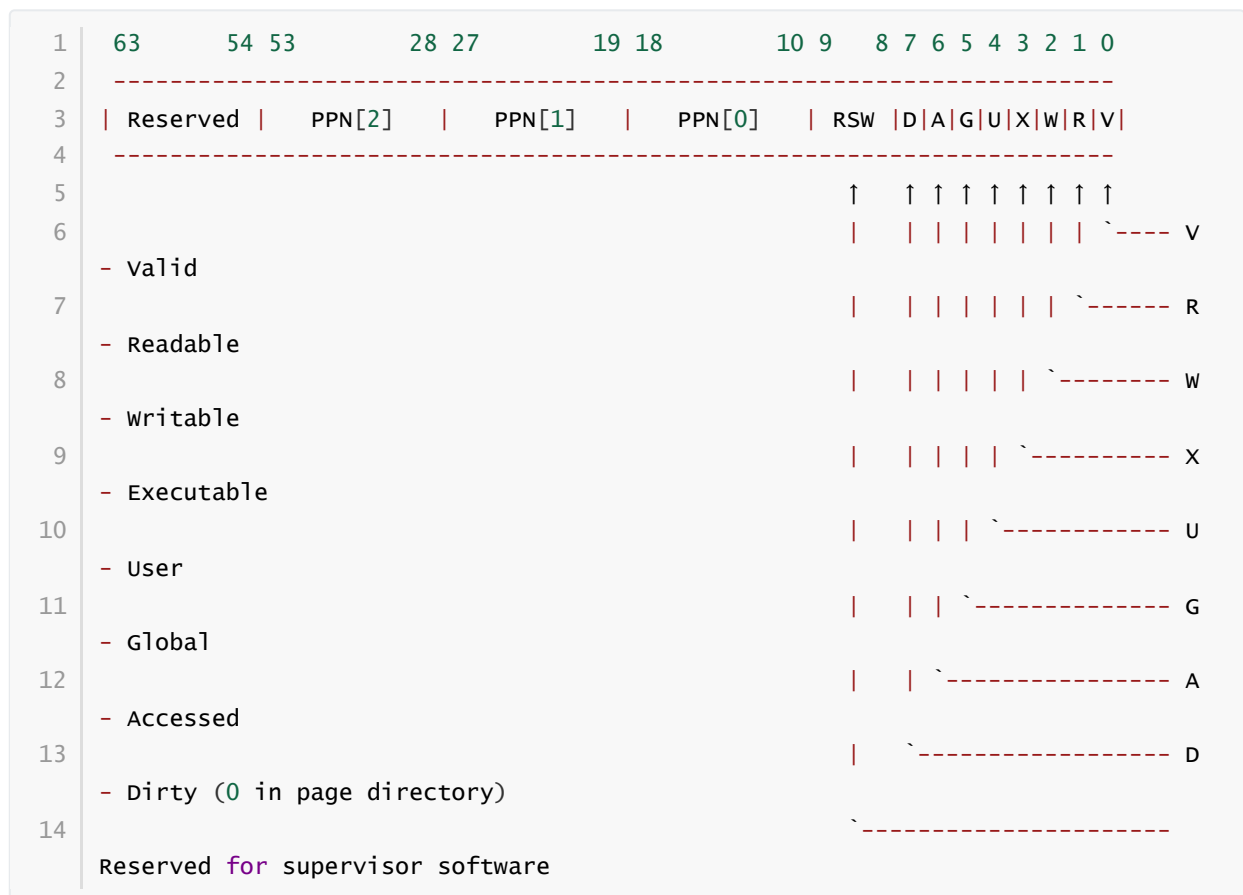


- Sv39 模式定义物理地址有 56 位，虚拟地址有 64 位。但是，虚拟地址的 64 位只有低 39 位有效。通过虚拟内存布局图我们可以发现，其 63-39 位为 0 时代表 user space address，为 1 时代表 kernel space address。
- Sv39 支持三级页表结构，VPN[2] VPN[1] VPN[0] (Virtual Page Number) 分别代表每级页表的虚拟页号，PPN[2] PPN[1] PPN[0] (Physical Page Number) 分别代表每级页表的物理页号。物理地址和虚拟地址的低12位表示页内偏移（page offset）。
- 具体介绍请阅读 [RISC-V Privileged Spec 4.4.1](#)。

3.3.3 RISC-V Sv39 Page Table Entry

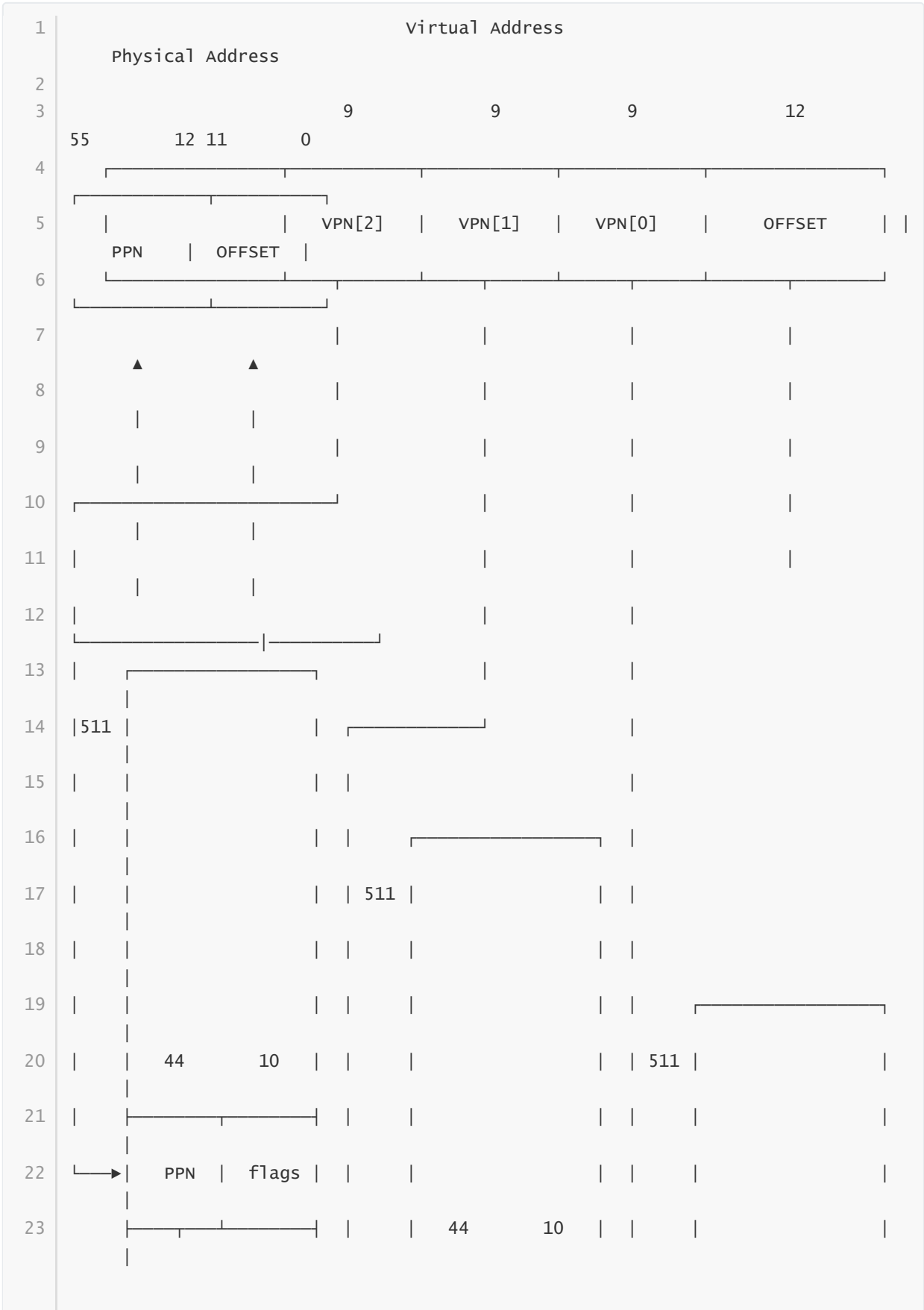


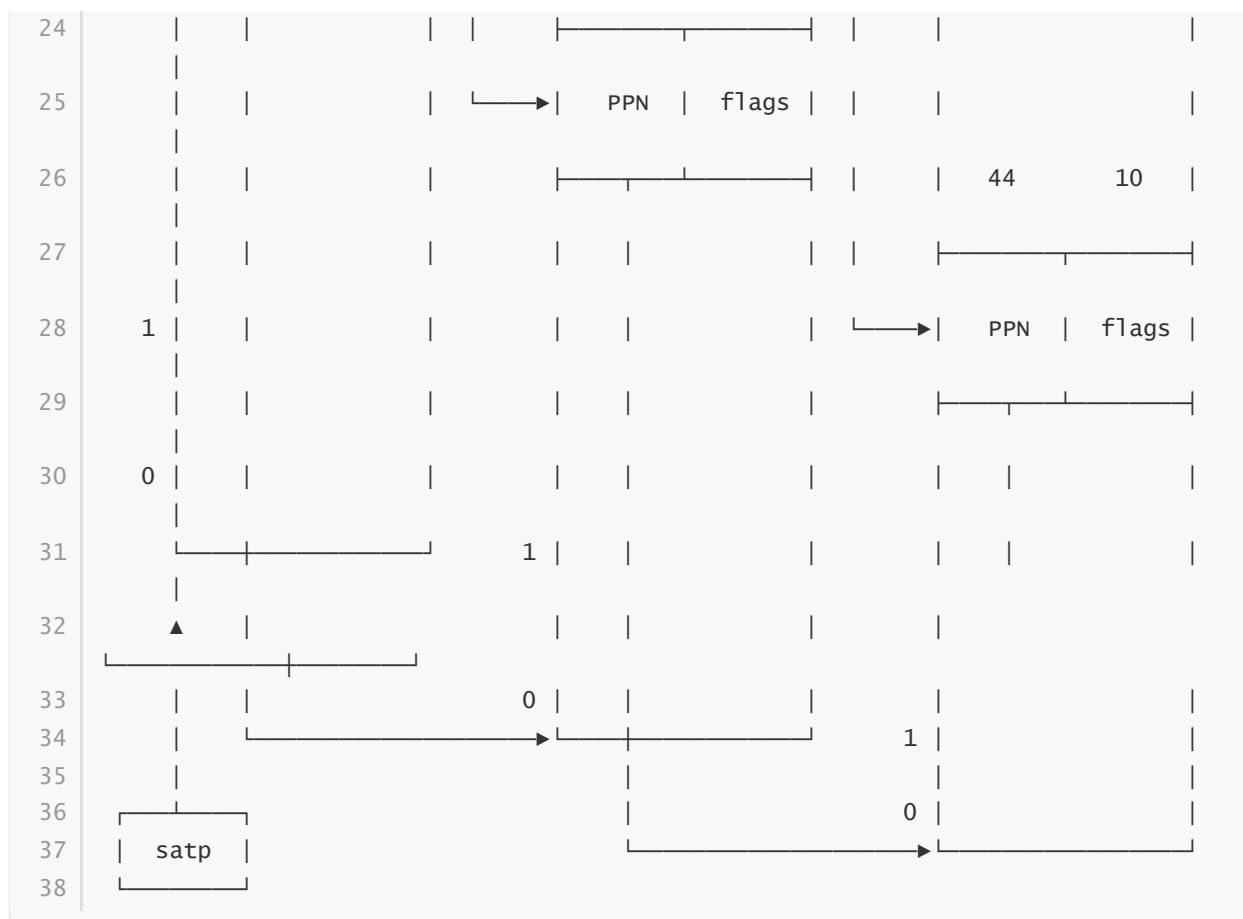
- 0 ~ 9 bit: protection bits
 - V: 有效位，当 V = 0，访问该 PTE 会产生 Pagefault。
 - R: R = 1 该页可读。
 - W: W = 1 该页可写。
 - X: X = 1 该页可执行。
 - U, G, A, D, RSW 本次实验中设置为 0 即可。

- 具体介绍请阅读 [RISC-V Privileged Spec 4.4.1](#)

3.3.4 RISC-V Address Translation

虚拟地址转化为物理地址流程图如下，具体描述见 [RISC-V Privileged Spec 4.3.2](#)：





4. 实验步骤

4.1 准备工程

- 修改 `defs.h`，在 `defs.h` 添加如下内容：

```

1 #define OPENSBI_SIZE (0x200000)
2
3 #define VM_START (0xffffffe000000000)
4 #define VM_END   (0xfffffffff00000000)
5 #define VM_SIZE   (VM_END - VM_START)
6
7 #define PA2VA_OFFSET (VM_START - PHY_START)

```

- 从 `repo` 同步以下代码: `vmlinux.lds`。并按照以下步骤将这些文件正确放置。

```

1 .
2 └─ arch
3     └─ riscv
4         └─ kernel
5             └─ vmlinux.lds

```

- 更改项目顶层目录的 `Makefile`，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，修改后的顶层目录的 `Makefile` 如下：

```

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) -mmodel=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     @echo -e '\n'Build Finished OK
21
22 TEST:
23     ${MAKE} -C lib all
24     ${MAKE} -C test test
25     ${MAKE} -C init all
26     ${MAKE} -C arch/riscv all
27     @echo -e '\n'Build Finished OK
28
29 run: all
30     @echo Launch the qemu .....
31     @qemu-system-riscv64 -nographic -machine virt -kernel vmlinux -bios
    default
32
33 test-run: TEST
34     @echo Launch the qemu .....
35     @qemu-system-riscv64 -nographic -machine virt -kernel vmlinux -bios
    default
36
37 debug: all
38     @echo Launch the qemu for debug .....
39     @qemu-system-riscv64 -nographic -machine virt -kernel vmlinux -bios
    default -S -s
40
41 test-debug: TEST
42     @echo Launch the qemu for debug .....
43     @qemu-system-riscv64 -nographic -machine virt -kernel vmlinux -bios
    default -S -s
44
45 clean:
46     ${MAKE} -C lib clean

```

```

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

```

4.2 开启虚拟内存映射。

4.2.1 setup_vm 的实现

- 将 0x80000000 开始的 1GB 区域进行两次映射，其中一次是等值映射 ($PA == VA$)，另一次是将其映射到 direct mapping area (使得 $PA + PV2VA_OFFSET == VA$)。
 - 将 PHY_START 右移30位，结果和 0x1ff 相与得到中间9位，作为 early_pgtbl 数组的索引
 - 将 early_pgtbl[index] 等值映射，对齐 PAGE_SIZE 后再和 15 按位或将权限 V | R | W | X 位设置为 1
 - 同样将 $PA + PV2VA_OFFSET$ 做映射

代码如下图所示：

```

1  /* early_pgtbl: 用于 setup_vm 进行 1GB 的映射。 */
2  unsigned long early_pgtbl[512] __attribute__((__aligned__(0x1000)));
3  void setup_vm(void) {
4      /*
5       1. 由于是进行 1GB 的映射 这里不需要使用多级页表
6       2. 将 va 的 64bit 作为如下划分： | high bit | 9 bit | 30 bit |
7          high bit 可以忽略dfd
8          中间9 bit 作为 early_pgtbl 的 index
9          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。
10         3. Page Table Entry 的权限 V | R | W | X 位设置为 1
11      */
12      uint64 index = PHY_START >> 30;
13      index = index & 0x1ff;
14      early_pgtbl[index] = PHY_START >> 2 | 15;
15
16      index = (PHY_START + PA2VA_OFFSET) >> 30;
17      index = index & 0x1ff;
18      early_pgtbl[index] = PHY_START >> 2 | 15;
19  }

```

- 完成上述映射之后，通过 relocate 函数，完成对 satp 的设置，以及跳转到对应的虚拟地址。
 - 首先将 ra 和 sp 的值都增加了偏移量 PA2VA_OFFSET，进入虚拟地址运行
 - 根据 pagetable 设置 satp 寄存器，减去偏移量使用物理地址，mode 位设为 8；PPN 位设为 $PA \gg 12$
 - flush tlb 和 icache

```

1  relocate:

```

```

2      # set ra = ra + PA2VA_OFFSET
3      # set sp = sp + PA2VA_OFFSET (If you have set the sp before)
4
5      #####
6      #   YOUR CODE HERE   #
7      li t0, 0xfffffffff8000000
8      add ra, ra, t0
9      add sp, sp, t0
10     #####
11
12     # set satp with early_pgtbl
13
14     #####
15     #   YOUR CODE HERE   #
16     la t0, early_pgtbl
17     li t1, 0xfffffffff8000000
18     sub t0, t0, t1
19     # PA >> 12 == PPN
20     addi t1, x0, 12
21     srl t0, t0, t1
22     # mode = 8
23     addi t2, x0, 1
24     slli t2, t2, 63
25     or t0, t0, t2
26     csrw satp, t0
27     #####
28
29     # flush tlb
30     sfence.vma zero, zero
31
32     # flush icache
33     fence.i
34
35     ret

```

4.2.2 setup_vm_final 的实现

- 由于 `setup_vm_final` 中需要申请页面的接口，应该在其之前完成内存管理初始化，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需要调整为虚拟地址，在原来的 `PHY_END` 加上偏移量 `PA2VA_OFFSET`

```

1 void mm_init(void) {
2     kfreerange(&_amp;kernel, (char *) (PHY_END+PA2VA_OFFSET));
3     printk("...mm_init done!\n");
4 }

```

- 对所有物理内存 (128M) 进行映射，并设置正确的权限，采用三级页表映射。

- 首先获取 text 段的起始物理地址、虚拟地址以及地址空间大小，使用下一段的地址 `_srodata` - `_stext` 计算地址空间大小，将 perm 设为 1011，即 X|-|R|V
- 调用 `create_mapping` 进行地址映射
- 与 text 段类似，依次计算 rodata 段和其他部分的起始地址和空间大小，rodata 的 perm 设为 0011，即 -|-|R|V，其他段的 perm 设为 0111，即 -|W|R|V，调用 `create_mapping` 进行地址映射

```

1 char _stext[];
2 char _srodata[];
3 char _sdata[];
4 char _sbss[];
5 char _kernel[];
6 void setup_vm_final(void) {
7     // mapping kernel text X|-|R|V
8     uint64 va_text = (uint64)_stext;
9     uint64 pa_text = (uint64)_stext - PA2VA_OFFSET;
10    uint64 sz_text = (uint64)_srodata - (uint64)_stext;
11    int perm_text = 0b1011;
12    create_mapping(swapper_pg_dir, va_text, pa_text, sz_text, perm_text);
13
14    // mapping kernel rodata -|-|R|V
15    uint64 va_rodata = (uint64)_srodata;
16    uint64 pa_rodata = (uint64)_srodata - PA2VA_OFFSET;
17    uint64 sz_rodata = (uint64)_sdata - (uint64)_srodata;
18    int perm_rodata = 0b0011;
19    create_mapping(swapper_pg_dir, va_rodata, pa_rodata, sz_rodata,
20    perm_rodata);
21
22    // mapping other memory -|W|R|V
23    uint64 va_orther = (uint64)_sdata;
24    uint64 pa_orther = (uint64)_sdata - PA2VA_OFFSET;
25    uint64 sz_orther = PHY_SIZE - OPENSBI_SIZE - ((uint64)_sdata -
26    (uint64)_stext);
27    int perm_orther = 0b0111;
28    create_mapping(swapper_pg_dir, va_orther, pa_orther, sz_orther,
29    perm_orther);
30    return;
31 }

```

- 对每段地址进行映射：
 - 通过 for 循环对该段虚拟地址进行映射，每次增加一页（一个 `PGSIZE`）。
 - 计算多级页表中的各级偏移量
 - 首先检查最高级的页表项是否存在，如果不存在（V bit为0），则使用 `ka1loc()` 分配一个新的页面作为下一级页表，并设置V bit为1。
 - 获取下一级页表的地址，将地址对齐，根据下一级的偏移量查找对应的地址，重复直到找到最终真正的地址
 - 进行下一页的映射

```

1  /* 创建多级页表映射关系 */
2  create_mapping(uint64 *pgtbl, uint64 va, uint64 pa, uint64 sz, int 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
7          perm 为映射的读写权限
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
13     for(uint64 i = va ; i < va + sz ; i += PGSIZE){
14         uint64 offset[3];
15         uint64 page[3];
16         offset[0] = va >> 12 & 0x1FF;
17         offset[1] = va >> 21 & 0x1FF;
18         offset[2] = va >> 30 & 0x1FF;
19         if(!pgtbl[offset[2]] & 1){
20             uint64 new_page_addr = kalloc();
21             new_page_addr -= 0xfffffffdf80000000;
22             printk(new_page_addr);
23             pgtbl[offset[2]] = (new_page_addr & 111111111111000) >> 2;
24             pgtbl[offset[2]] |= 0x0000000000000001;
25         }
26         page[2] = pgtbl[offset[2]];
27
28         uint64* pgtbl_next = (uint64*)((page[2] & 111111111111000) << 2);
29         if(!pgtbl[offset[1]] & 1){
30             uint64 new_page_addr = kalloc();
31             new_page_addr -= 0xfffffffdf80000000;
32             pgtbl[offset[1]] = (new_page_addr & 111111111111000) >> 2;
33             pgtbl[offset[1]] |= 0x0000000000000001;
34         }
35         page[1] = pgtbl[offset[1]];
36
37         uint64* pgtbl_nnext = (uint64*)((page[1] & 111111111111000) << 2);
38         page[0] = ((pa & 111111111111000) << 2) | perm;
39         pgtbl_nnext[offset[0]] = page[0];
40         va += PGSIZE;
41         pa += PGSIZE;
42     }
43     return;
44 }

```

- 在 head.S 中 适当的位置调用 `setup_vm_final`、`setup_vm_final`、`setup_vm`、`relocate`
将栈顶指针减去偏移量获得物理地址，后续调用 `relocate` 再转化为虚拟地址

```
1  _start:
2      #将栈顶指针放入sp
3      la sp, boot_stack_top
4      li t0, 0xffffffffdf80000000
5      sub sp, sp, t0
6      call setup_vm
7      call relocate
8      call mm_init
9      jal setup_vm_final
10     jal task_init
11     jal test_init
```

4.3 编译及测试

- 运行后的结果如下图所示，与lab2相同

[illegible]

思考题

1. 验证 `.text`、`.rodata` 段的属性是否成功设置，给出截图。
 - 在GDB中调试程序，先在 `start_kernel` 处打上断点，再运行到 `start_kernel`

```
(gdb) b start_kernel
Breakpoint 1 at 0xffffffff0002011c8: file main.c, line 5.
(gdb) continue
Continuing.

Breakpoint 1, start_kernel () at main.c:5
5      printk("2023");
(gdb) █
```

- 使用 `info files` 命令查看文件信息：

```
(gdb) info files
Symbols from "/mnt/d/study/OS/os23fall-stu/src/lab3/vmlinux".
Remote target using gdb-specific protocol:
  ~/mnt/d/study/OS/os23fall-stu/src/lab3/vmlinux', file type elf64-littleriscv.
Entry point: 0xffffffff000200000
0xffffffff000200000 - 0xffffffff000202598 is .text
0xffffffff000203000 - 0xffffffff00020334e is .rodata
0xffffffff000204000 - 0xffffffff000204008 is .data
0xffffffff000204008 - 0xffffffff000204058 is .got
0xffffffff000204058 - 0xffffffff000204068 is .got.plt
0xffffffff000205000 - 0xffffffff00020a1c5 is .bss
While running this, GDB does not access memory from...
Local exec file:
  ~/mnt/d/study/OS/os23fall-stu/src/lab3/vmlinux', file type elf64-littleriscv.
Entry point: 0xffffffff000200000
0xffffffff000200000 - 0xffffffff000202598 is .text
0xffffffff000203000 - 0xffffffff00020334e is .rodata
0xffffffff000204000 - 0xffffffff000204008 is .data
0xffffffff000204008 - 0xffffffff000204058 is .got
0xffffffff000204058 - 0xffffffff000204068 is .got.plt
0xffffffff000205000 - 0xffffffff00020a1c5 is .bss
```

可以看到 `.text` , `.rodata` 的地址为虚拟地址

- 使用 `maintenance info sections` 命令来查看详细的段信息

```
(gdb) maintenance info sections
Exec file: ~/mnt/d/study/OS/os23fall-stu/src/lab3/vmlinux', file type elf64-littleriscv.
[0] 0xffffffff000200000->0xffffffff000202598 at 0x00001000: .text ALLOC LOAD READONLY CODE HAS_CONTENTS
[1] 0xffffffff000203000->0xffffffff00020334e at 0x00004000: .rodata ALLOC LOAD READONLY DATA HAS_CONTENTS
[2] 0xffffffff000204000->0xffffffff000204008 at 0x00005000: .data ALLOC LOAD DATA HAS_CONTENTS
[3] 0xffffffff000204008->0xffffffff000204058 at 0x00005008: .got ALLOC LOAD DATA HAS_CONTENTS
[4] 0xffffffff000204058->0xffffffff000204068 at 0x00005058: .got.plt ALLOC LOAD DATA HAS_CONTENTS
[5] 0xffffffff000205000->0xffffffff00020a1c5 at 0x00005068: .bss ALLOC
[6] 0x00000000->0x0000018fb at 0x00005068: .debug_info READONLY HAS_CONTENTS
[7] 0x00000000->0x00000d1a at 0x00006963: .debug_abbrev READONLY HAS_CONTENTS
[8] 0x00000000->0x000003c0 at 0x00007680: .debug_aranges READONLY HAS_CONTENTS
[9] 0x00000000->0x00000285 at 0x00007a40: .debug_rnglists READONLY HAS_CONTENTS
[10] 0x00000000->0x00000209f at 0x00007cc5: .debug_line READONLY HAS_CONTENTS
[11] 0x00000000->0x000000549 at 0x00009d64: .debug_str READONLY HAS_CONTENTS
[12] 0x00000000->0x000001b5 at 0x0000a2ad: .debug_line_str READONLY HAS_CONTENTS
[13] 0x00000000->0x0000002b at 0x0000a462: .comment READONLY HAS_CONTENTS
[14] 0x00000000->0x000000032 at 0x0000a48d: .riscv.attributes READONLY HAS_CONTENTS
[15] 0x00000000->0x000000660 at 0x0000a4c0: .debug_frame READONLY HAS_CONTENTS
```

可以看到：

`.text` 的属性为 `ALLOC` (段在进程的虚拟地址空间中有分配的空间)；`LOAD` (段从文件加载到内存中)；`READONLY` (段是只读的)；`CODE` (段包含可执行代码)

`.rodata` 的属性为 `ALLOC` (段在进程的虚拟地址空间中有分配的空间)；`LOAD` (段从文件加载到内存中)；`READONLY` (段是只读的)；`DATA` (段包含只读数据)

属性设置正确

2. 为什么我们在 `setup_vm` 中需要做等值映射？

答：在建立三级页表的时候需要通过物理页号在页表中找到物理地址，此时的运行在虚拟地址上，直接使用物理地址会产生访问错误，所以需要进行等值映射

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。

答：在通过物理页号得到下一级页表中的物理地址时直接加上偏移量 `PA2VA_OFFSET`，得到对应的虚拟地址，就不用在 `setup_vm` 中做等值映射了

```
1 uint64* pgtbl_next = (uint64*)((page[2] & 111111111111000) << 2 +  
  PA2VA_OFFSET);  
2 uint64* pgtbl_nnext = (uint64*)((page[1] & 111111111111000) << 2 +  
  PA2VA_OFFSET);
```