

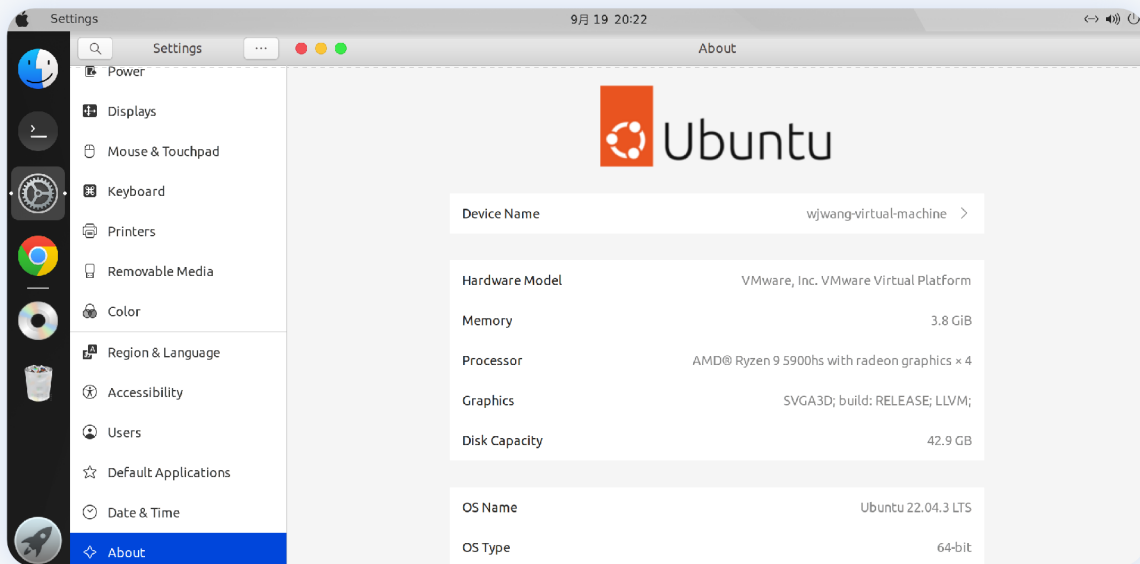
# Lab 5: RV64 缺页异常处理

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lab4中我已经处理好了 `thread_info` 的问题, 在此无需调整

## 实现 VMA

修改 `proc.h` , 增加如下相关结构

```
1  #define VM_X_MASK          0x0000000000000008
2  #define VM_W_MASK          0x0000000000000004
3  #define VM_R_MASK          0x0000000000000002
4  #define VM_ANONYM          0x0000000000000001
5
6  struct vm_area_struct {
7      uint64_t vm_start;      /* VMA 对应的用户态虚拟地址的开始 */
8      uint64_t vm_end;        /* VMA 对应的用户态虚拟地址的结束 */
9      uint64_t vm_flags;      /* VMA 对应的 flags */
10
11     /* uint64_t file_offset_on_disk */ /* 原本需要记录对应的文件在磁盘上的位置,
12                                         但是我们只有一个文件 uapp, 所以暂时不需要记录 */
13
14     uint64_t vm_content_offset_in_file; /* 如果对应了一个文件,
15                                         那么这块 VMA 起始地址对应的文件内容相对文件起始位置的偏移量,
16                                         也就是 ELF 中各段的 p_offset 值 */
17
18     uint64_t vm_content_size_in_file; /* 对应的文件内容的长度。
19                                         思考为什么还需要这个域?
20                                         和 (vm_end-vm_start)
21                                         一比, 不是冗余了吗? */
22 };
23
24 struct task_struct {
25     uint64_t state;
26     uint64_t counter;
27     uint64_t priority;
28     uint64_t pid;
29
30     struct thread_struct thread;
31     pagetable_t pgd;
32
33     uint64_t vma_cnt; /* 下面这个数组里的元素的数量 */
34     struct vm_area_struct vmas[0]; /* 为什么可以开大小为 0 的数组?
35                                     这个定义可以和前面的 vma_cnt 换个位置吗? */
36 };
```

为了支持 `Demand Paging` , 我们支持对 `vm_area_struct` 的添加和查找。

```

2  /* 创建一个新的 vma */
3  void do_mmap(struct task_struct *task, uint64_t addr, uint64_t length,
4             uint64_t flags,
5             uint64_t vm_content_offset_in_file, uint64_t
6             vm_content_size_in_file) {
7      struct vm_area_struct *vma = &(task->vmas[task->vma_cnt++]); // find last
8      vma
9      vma->vm_start = addr;
10     vma->vm_end = addr + length;
11     vma->vm_flags = flags;
12     vma->vm_content_offset_in_file = vm_content_offset_in_file;
13     vma->vm_content_size_in_file = vm_content_size_in_file;
14 }
15
16 /* 查找包含某个 addr 的 vma, 该函数主要在 Page Fault 处理时起作用 */
17 struct vm_area_struct *find_vma(struct task_struct *task, uint64_t addr) {
18     for (int i = 0; i < task->vma_cnt; ++i) {
19         if (task->vmas[i].vm_start ≤ addr && addr < task->vmas[i].vm_end) {
20             return &(task->vmas[i]);
21         }
22     }
23     return NULL;
24 }

```

## Demand Paging

在初始化一个 task 时我们既不分配内存，又不更改页表项来建立映射。回退到用户态进行程序执行的时候就会因为没有映射而发生 Page Fault，进入我们的 Page Fault Handler 后，我们再分配空间（按需要拷贝内容）进行映射。

我们需要修改 `task_init` 中的 `load_program` 函数代码，更改为 Demand Paging，将两次 `create_mapping` 改为 `do_mmap`，取消分配内存的操作：

```

1  static uint64_t load_program(struct task_struct* task) {
2      Elf64_Ehdr* ehdr = (Elf64_Ehdr*)_sramdisk;
3
4      if (ehdr->e_ident[0] ≠ 0x7f || ehdr->e_ident[1] ≠ 'E' ||
5          ehdr->e_ident[2] ≠ 'L' || ehdr->e_ident[3] ≠ 'F') {
6          printk("not elf format\n");
7          return -1;
8      }

```

```

9
10     uint64_t phdr_start = (uint64_t)ehdr + ehdr->e_phoff;
11     int phdr_cnt = ehdr->e_phnum;
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             // do mapping
20             uint64_t flags = (phdr->p_flags & PF_R) >> 1 | (phdr->p_flags &
PF_W) << 1 | (phdr->p_flags & PF_X) << 3;
21             // char* va = (char*)alloc_pages((size + PGSIZE - 1) / PGSIZE);
22             // // memcpy(va + offset, (char*)ehdr + phdr->p_offset, phdr-
>p_filesz);
23             // for (int j = 0; j < phdr->p_filesz; ++j) {
24             //     va[offset + j] = ((char*)ehdr)[phdr->p_offset + j];
25             // }
26             // memset(va + offset + phdr->p_filesz, 0, phdr->p_memsz - phdr-
>p_filesz);
27             do_mmap(task, phdr->p_vaddr, phdr->p_memsz, flags, phdr-
>p_offset, phdr->p_filesz);
28             // create_mapping(task->pgd, phdr->p_vaddr, (uint64_t)va -
PA2VA_OFFSET, size, 1 << 4 | 1 << 0 | (phdr->p_flags & 0x4) >> 1 | (phdr-
>p_flags & 0x2) << 1 | (phdr->p_flags & 0x1) << 3);
29         }
30     }
31
32     // allocate user stack and do mapping
33
34     // uint64_t user_stack = (uint64_t)alloc_page();
35     do_mmap(task, USER_END - PGSIZE, PGSIZE, 0x7, 0, 0);
36     // create_mapping(task->pgd, USER_END - PGSIZE, user_stack -
PA2VA_OFFSET, PGSIZE, 0x17);
37     // pc for the user program
38     task->thread.sepc = ehdr->e_entry;
39     return 0;
40 }

```

修改 `trap_handler`，添加跳转 `page_fault` 函数的部分，并在必要的时候进入 `unhandled_trap`：

```

1 void trap_handler(unsigned long scause, unsigned long sepc, struct pt_regs
  *regs) {
2     unsigned long temp = 1;
3     if(scause & (temp << 63)) { // interrupt
4         switch (scause & ~(temp << 63))
5         {
6             case 0x5:
7                 printk("[INTERRUPT] S mode timer interrupt!\n");
8                 clock_set_next_event();
9                 do_timer();
10                break;
11
12            default:
13                unhandled_trap(regs);
14                break;
15        }
16    } else {
17        switch (scause & ~(temp << 63))
18        {
19            case 0x8: // ECALL_FROM_U_MODE
20                uint64 syscall_id = regs->a7;
21                switch (syscall_id)
22                {
23                    case SYS_WRITE: // sys_write
24                        regs->a0 = sys_write(regs->a0, regs->a1, regs->a2);
25                        break;
26                    case SYS_GETPID: // sys_getpid
27                        regs->a0 = sys_getpid();
28                        break;
29                }
30                regs->sepc += 4;
31                break;
32
33            case 0xc:
34            case 0xd:
35            case 0xf:
36                do_page_fault(regs);
37                break;
38
39            default:

```

```

40         printk("Unhandled exception in %lx!\n", scause);
41         unhandled_trap(regs);
42         break;
43     }
44 } // exception
45 }

```

`unhandled_trap` 定义了发生未知错误时的处理方式：暂停等待：

```

1 void unhandled_trap(struct pt_regs *regs) {
2     printk("Unhandled trap:\n");
3     printk("sepc: %lx, scause: %lx, stval: %lx\n", regs->sepc, regs->scause,
4         regs->stval);
5     while(1);
6 }

```

进入 `page_fault` 之后处理步骤如下：

- 通过 `stval` 获得访问出错的虚拟内存地址（Bad Address）
- 通过 `find_vma()` 查找 Bad Address 是否在某个 `vma` 中
  - 找不到这个 `vma`，或者检查是否是权限错误，如果是这样我们不进行处理，直接进入 `unhandled_trap` .STOP
  - 找到了对应 `vma`，继续下一步
- 分配一个页，将这个页映射到对应的用户地址空间，并将这个页全部清空（之后不需要进行清零处理）
- 通过 ( `vma->vm_flags & VM_ANONYM` ) 获得当前的 VMA 是否是匿名空间
  - 如果 VMA 匿名，由于上一步已经进行清零操作，直接返回 .STOP
  - 如果不是匿名，并且拷贝页面存在 ( `end_seg > start_bad_seg_page` ) 我们将 `bad_addr` 所在的那一页拷贝过来，拷贝数据的开头有两种情况：
    - `start_seg` 和 `bad_addr` 在同一页，这时拷贝数据的开头不需要是整页的起始地址，而是需要从 `start_seg` 开始拷贝，即算出对应 `offset=start_seg`后12位
    - `start_seg` 和 `bad_addr` 不在同一页，这时直接从 `bad_addr` 对应文件中的那一页起始地址，从 `start_bad_seg_page` 开始拷贝，即置 `offset` 为零

结尾有以下两种情况（其中下方参考图的①号位置为 `start_bad_seg_page + PGSIZE`）：

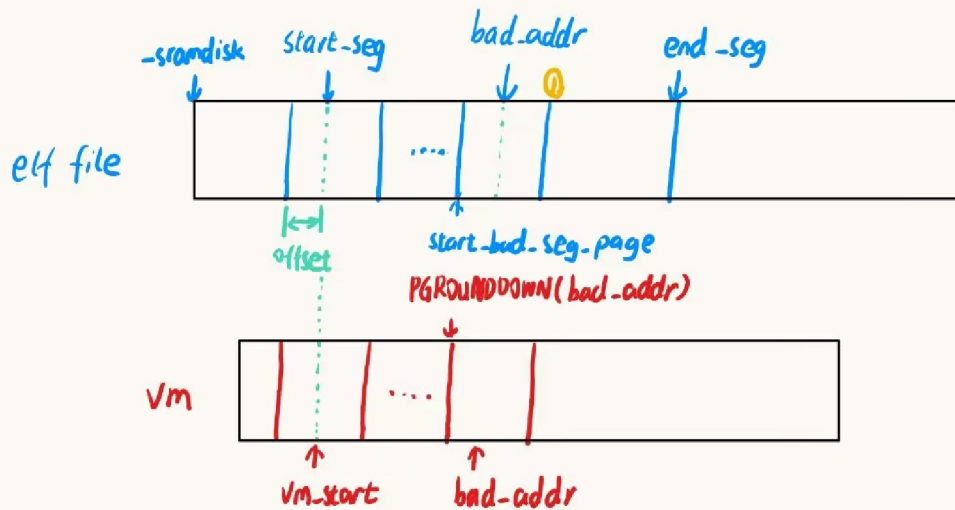
- `end_seg > start_bad_seg_page + PGSIZE`，这时我们直接拷贝 [ `start_bad_seg_page + offset`, `start_bad_seg_page + offset + PGSIZE` ]

- `start_bad_seg_page < end_seg < start_bad_seg_page + PGSIZE` , 这时我们拷贝[[ `start_bad_seg_page + offset` , `end_seg` ]

.STOP

- 如果不是匿名, 并且拷贝页面不存在 ( `end_seg ≤ start_bad_seg_page` ), 则不需要做处理, 上一步已经进行清零操作 .STOP

对于找到对应vma、 `start_seg` 和 `bad_addr` 不在同一页、 `start_bad_seg_page < end_seg < start_bad_seg_page + PGSIZE` 的情况, 我绘制了拷贝参考图进行辅助说明:



```

1 void do_page_fault(struct pt_regs *regs) {
2     uint64 bad_addr = regs->stval;
3     struct vm_area_struct *vma = find_vma(current, bad_addr);
4     // 范围检查 权限检查
5     if(vma == NULL ||
6         (regs->scause == 0xc && !(vma->vm_flags & VM_X_MASK)) ||
7         (regs->scause == 0xd && !(vma->vm_flags & VM_R_MASK)) ||
8         (regs->scause == 0xf && !(vma->vm_flags & VM_W_MASK))) {
9         printk("Bad Address: %lx\n", bad_addr);
10        unhandled_trap(regs);
11    }
12    uint64 page = alloc_page();
13    create_mapping(current->pgd, bad_addr, page - PA2VA_OFFSET, PGSIZE, (vma->vm_flags & 0b1110) | 0x11); // UXWRV
14    memset((void *)page, 0, PGSIZE);
15    if (!(vma->vm_flags & VM_ANONYM)) {

```

```

16         // segment start location in file
17         uint64 start_seg = (uint64)(_sramdisk) + vma-
>vm_content_offset_in_file;
18         uint64 offset = 0;
19         uint64 end_seg= (uint64)(_sramdisk) + vma->vm_content_offset_in_file
+vma->vm_content_size_in_file;
20         uint64 start_bad_seg_page = start_seg + PGROUNDDOWN(bad_addr) - vma-
>vm_start;
21         uint64 start_bad_seg = start_bad_seg_page;
22         if (PGROUNDDOWN(start_seg) == PGROUNDDOWN(bad_addr)) {
23             start_bad_seg = start_seg;
24             offset = start_seg & 0xfff;
25         }
26         if (end_seg ≤ start_bad_seg_page + PGSIZE && end_seg >
start_bad_seg_page){
27             uint64 valid_seg = end_seg - start_bad_seg_page - offset;
28             memcpy((void *)(page + offset), (void *)(start_bad_seg),
valid_seg);
29         } else if (end_seg > start_bad_seg_page + PGSIZE) {
30             memcpy((void *)(page + offset), (void *)(start_bad_seg), PGSIZE -
offset);
31         }
32     }
33 }

```

## 编译及测试

进程第一次被调度，`main` 作为 `uapp` 的情况下，一共会发生3次 Page Fault：



```

...proc_init done!
2022 Hello RISC-V
sstatus = 8000000000000000
sscratch = 0

switch to [PID = 1 COUNTER = 4 PRIORITY = 37]
Page Fault: sepc: 00000000000100e8, scause: 000000000000000c, stval: 00000000000100e8
Page Fault: sepc: 0000000000010124, scause: 000000000000000f, stval: 0000003fffffffff8
Page Fault: sepc: 0000000000010140, scause: 000000000000000d, stval: 0000000000011880
[PID = 1] is running, variable: 0
[PID = 1] is running, variable: 1
[INTERRUPT] S mode timer interrupt!
[PID = 1] is running, variable: 2
[PID = 1] is running, variable: 3
[INTERRUPT] S mode timer interrupt!
[PID = 1] is running, variable: 4
[PID = 1] is running, variable: 5
[INTERRUPT] S mode timer interrupt!
[PID = 1] is running, variable: 6
[PID = 1] is running, variable: 7
[PID = 1] is running, variable: 8
[INTERRUPT] S mode timer interrupt!

switch to [PID = 3 COUNTER = 8 PRIORITY = 52]
Page Fault: sepc: 00000000000100e8, scause: 000000000000000c, stval: 00000000000100e8
Page Fault: sepc: 0000000000010124, scause: 000000000000000f, stval: 0000003fffffffff8
Page Fault: sepc: 0000000000010140, scause: 000000000000000d, stval: 0000000000011880
[PID = 3] is running, variable: 0
[PID = 3] is running, variable: 1

```

第二次被调度，不会发生page fault:

```

[INTERRUPT] S mode timer interrupt!
SET [PID = 1 COUNTER = 1]
SET [PID = 2 COUNTER = 4]
SET [PID = 3 COUNTER = 10]

switch to [PID = 1 COUNTER = 1 PRIORITY = 37]
[PID = 1] is running, variable: 9
[PID = 1] is running, variable: 10
[INTERRUPT] S mode timer interrupt!

switch to [PID = 2 COUNTER = 4 PRIORITY = 88]
[PID = 2] is running, variable: 19
[PID = 2] is running, variable: 20
[INTERRUPT] S mode timer interrupt!
[PID = 2] is running, variable: 21
[PID = 2] is running, variable: 22
[INTERRUPT] S mode timer interrupt!
[PID = 2] is running, variable: 23
[PID = 2] is running, variable: 24
[INTERRUPT] S mode timer interrupt!
[PID = 2] is running, variable: 25
[PID = 2] is running, variable: 26
[INTERRUPT] S mode timer interrupt!

switch to [PID = 3 COUNTER = 10 PRIORITY = 52]
[PID = 3] is running, variable: 17
[PID = 3] is running, variable: 18
[INTERRUPT] S mode timer interrupt!
[PID = 3] is running, variable: 19

```

## 思考题

1. `uint64_t vm_content_size_in_file;` 对应的文件内容的长度。为什么还需要这个域?

实验四中, 我们知道 `p_filesz < p_memsz`, 因为elf文件不会为一些特定段预留空间。`vm_start - vm_end` 相当于 `p_memsz`, 而 `vm_content_size_in_file` 相当于 `p_filesz`。我们使用 `vm_content_size_in_file` 来确定需要清零的区域, 内存中非匿名区域大于 `vm_content_size_in_file` 的部分需要清零。

2. `struct vm_area_struct vmas[0];` 为什么可以开大小为 0 的数组? 这个定义可以和前面的 `vma_cnt` 换个位置吗?

C语言中, 声明大小为0的数组是非标准的做法, 但是C99标准中引入柔性数组成员, 允许在结构体的最后一个成员中使用未知大小的数组, 最后一个零长数组实际上充当了一个指向结构体末尾连续内存块的指针, 可以根据需要动态地分配更多的内存给这个结构体, 以存储可变长度的数据。所以这个定义与前面的元素不可更换位置, 如果更换, 首先不能根据 `task_struct` 的地址直接找到 `vma_cnt` 以及 `vmas` 的位置, 也不符合C语言的标准。