

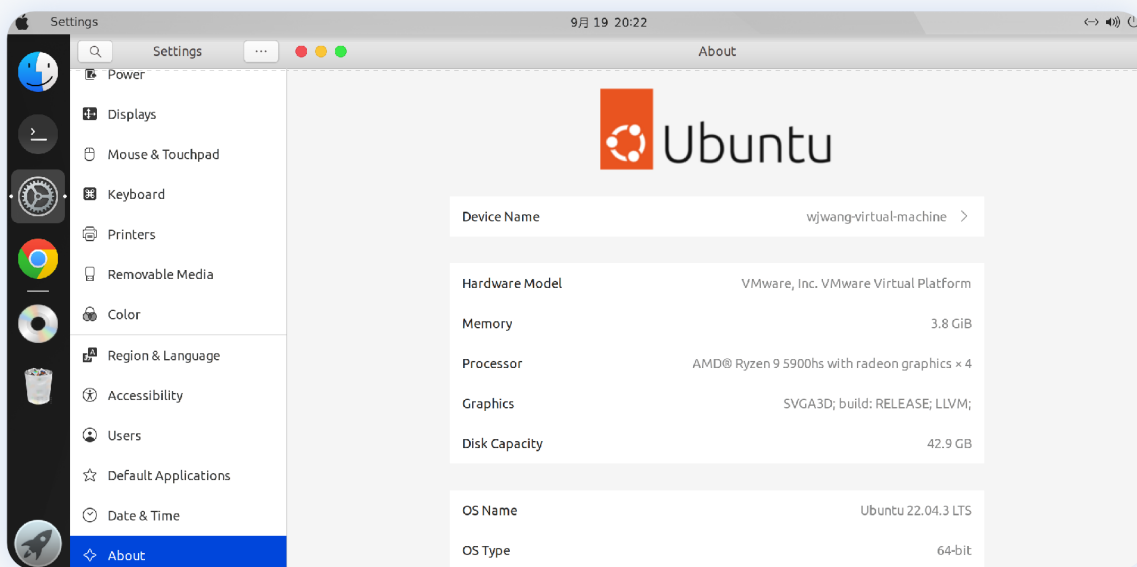
# Lab 6: 实现 fork 机制

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
```



## 调用sys\_clone

添加 `sys_clone` 的代码:

```
1  // syscall.h
2  #define SYS_CLONE 220
```

在 `trap_handler` 中添加处理 `sys_clone` 的代码:

```

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         // ...
5     } else {
6         switch (scause & ~(temp << 63))
7         {
8             case 0x8: // ECALL_FROM_U_MODE
9                 uint64 syscall_id = regs->a7;
10                switch (syscall_id)
11                {
12                    // ...
13                    case SYS_CLONE: // sys_clone
14                        regs->a0 = sys_clone(regs);
15                        break;
16                }
17                // ...
18            }
19        } // exception
20    }

```

## \_\_ret\_from\_fork

在 `_traps` 中的 `jal x1, trap_handler` 后面插入一个符号：

```

1     .global _traps
2     _traps:
3     ...
4     jal x1, trap_handler
5     .global __ret_from_fork
6     __ret_from_fork:
7     ... ;利用 sp 从栈中恢复出寄存器的值
8     sret

```

## 实现sys\_clone

`sys_clone` 的处理逻辑如下：

- 创建一个新的 task，将的 parent task 的整个页复制到新创建的 `task_struct` 页上，对复制过来的 `task_struct` 重新设置 `pid`，设置 `thread.ra` 设置为 `__ret_from_fork`
- 根据parent task的 `regs` 位置计算出child task对应的 `regs` 位置，将 `thread.sp` 设置为改地址
- 设置child task的 `pt_regs`，子进程的 `fork` 返回值 `a0` 设为0，`sepc` 设置为父进程 `sepc + 4`，使子进程跳过 `fork` 的这一条指令，继续运行下一条指令
- 为child task分配一个根页表，由于之前我们的 `task_init` 中的进程都是共享页表的，所以在这里，我也不创建内核页表的映射，直接使用之前的方式，让它们共享内核页表
- 最后，我遍历父进程的 `vma`，找到父进程对应 `vma` 中valid的PTE，将这些PTE对应的frame拷贝出一份新的，即child task在用户态会用到的内存

```

1  int IsValid(uint64 *pgd, uint64 va) {
2      uint64 VPN[3];
3      VPN[2] = (va >> 30) & 0x1ff; // 9 bit
4      VPN[1] = (va >> 21) & 0x1ff;
5      VPN[0] = (va >> 12) & 0x1ff;
6      uint64 *pte = pgd;
7      for (int j = 2; j > 0; j--) {
8          if ((pte[VPN[j]] & 0x1) == 0) {
9              return 0;
10         }
11         else {
12             pte = (uint64 *)((pte[VPN[j]] >> 10 << 12) + PA2VA_OFFSET);
13         }
14     }
15     if ((pte[VPN[0]] & 0x1) == 0) {
16         return 0;
17     }
18     return 1;
19 }
20
21 uint64_t sys_clone(struct pt_regs *regs) {
22     // 创建一个新的 task
23     struct task_struct *child = (struct task_struct *)kalloc();
24     memcpy((void*)child, (void*)current, PGSIZE);
25     for (int i = 0; i < NR_TASKS; ++i) {
26         if (task[i] == NULL) {
27             task[i] = child;
28             child->pid = i;
29             // printk("sys_clone: child->pid = %d\n", child->pid);
30             break;
31         }

```

```

32     }
33     child→thread.ra = (uint64)__ret_from_fork;
34
35     // 利用参数 regs 来计算出 child task 的对应的 pt_regs 的地址
36     uint64 offset = (uint64)regs - PGROUNDDOWN((uint64)regs);
37     struct pt_regs *child_regs = (struct pt_regs *)((uint64)child + offset);
38     child→thread.sp = (uint64)child_regs;
39     child_regs→a0 = 0;
40     // child_regs→sscratch = regs→sscratch;
41     child_regs→sepc = regs→sepc + 4;
42
43     // 为 child task 分配一个根页表
44     child→pgd = (pagetable_t)alloc_page();
45     memcpy((void *)child→pgd, (void *)swapper_pg_dir, PGSIZE);
46
47     // 根据 parent task 的页表和 vma 来分配并拷贝 child task 在用户态会用到的内存
48     for (int i = 0; i < current→vma_cnt; ++i) {
49         struct vm_area_struct *vma = &current→vmas[i];
50         uint64 now_page = PGROUNDDOWN(vma→vm_start);
51         while (now_page < vma→vm_end) {
52             if (IsValid(current→pgd, now_page)) {
53                 // printk("sys_clone: now_page = %lx\n", now_page);
54                 uint64 page = alloc_page();
55                 memcpy((void *)page, (void *)now_page, PGSIZE);
56                 create_mapping(child→pgd, now_page, page - PA2VA_OFFSET,
PGSIZE, (vma→vm_flags & 0b1110) | 0x11); // UXWRV
57             }
58             now_page += PGSIZE;
59         }
60     }
61
62     return child→pid;
63
64 }

```

## 编译及测试

第一个main函数:

```

...setup_vm done!
...buddy_init done!
...setup_vm_final done!
...proc_init done!
2022 Hello RISC-V
sstatus = 8000000000006000
sscratch = 0

switch to [PID = 1 COUNTER = 4 PRIORITY = 37]
Page Fault: sepc: 00000000000100e8, scause: 000000000000000c, stval: 00000000000100e8
Page Fault: sepc: 0000000000010158, scause: 000000000000000f, stval: 0000003fffffff8
[S] New task: 2
Page Fault: sepc: 00000000000101f0, scause: 000000000000000d, stval: 0000000000011978
[U-PARENT] pid: 1 is running!, global_variable: 0
[U-PARENT] pid: 1 is running!, global_variable: 1
[U-PARENT] pid: 1 is running!, global_variable: 2
[INTERRUPT] S mode timer interrupt!
[U-PARENT] pid: 1 is running!, global_variable: 3
[U-PARENT] pid: 1 is running!, global_variable: 4
[INTERRUPT] S mode timer interrupt!
[U-PARENT] pid: 1 is running!, global_variable: 5
[U-PARENT] pid: 1 is running!, global_variable: 6
[INTERRUPT] S mode timer interrupt!
[U-PARENT] pid: 1 is running!, global_variable: 7
[U-PARENT] pid: 1 is running!, global_variable: 8
[INTERRUPT] S mode timer interrupt!

switch to [PID = 2 COUNTER = 4 PRIORITY = 37]
Page Fault: sepc: 000000000001018c, scause: 000000000000000d, stval: 0000000000011978
[U-CHILD] pid: 2 is running!, global_variable: 0
[U-CHILD] pid: 2 is running!, global_variable: 1
[U-CHILD] pid: 2 is running!, global_variable: 2
[INTERRUPT] S mode timer interrupt!
[U-CHILD] pid: 2 is running!, global_variable: 3
[U-CHILD] pid: 2 is running!, global_variable: 4
[INTERRUPT] S mode timer interrupt!
[U-CHILD] pid: 2 is running!, global_variable: 5
[U-CHILD] pid: 2 is running!, global_variable: 6
[INTERRUPT] S mode timer interrupt!
[U-CHILD] pid: 2 is running!, global_variable: 7
[U-CHILD] pid: 2 is running!, global_variable: 8
[INTERRUPT] S mode timer interrupt!
SET [PID = 1 COUNTER = 1]
SET [PID = 2 COUNTER = 4]

switch to [PID = 1 COUNTER = 1 PRIORITY = 37]
[U-PARENT] pid: 1 is running!, global_variable: 9
[U-PARENT] pid: 1 is running!, global_variable: 10

```

第二个main函数:

```

2022 Hello RISC-V
ssstatus = 8000000000006000
sscratch = 0

switch to [PID = 1 COUNTER = 4 PRIORITY = 37]
Page Fault: sepc: 0000000000100e8, scause: 00000000000000c, stval: 0000000000100e8
Page Fault: sepc: 000000000010158, scause: 00000000000000f, stval: 0000003fffffff8
Page Fault: sepc: 00000000001017c, scause: 00000000000000d, stval: 000000000011a00
[U] pid: 1 is running!, global_variable: 0
[U] pid: 1 is running!, global_variable: 1
[U] pid: 1 is running!, global_variable: 2
[S] New task: 2
[U-PARENT] pid: 1 is running!, global_variable: 3
[U-PARENT] pid: 1 is running!, global_variable: 4
[U-PARENT] pid: 1 is running!, global_variable: 5
[INTERRUPT] S mode timer interrupt!
[U-PARENT] pid: 1 is running!, global_variable: 6
[U-PARENT] pid: 1 is running!, global_variable: 7
[INTERRUPT] S mode timer interrupt!
[U-PARENT] pid: 1 is running!, global_variable: 8
[U-PARENT] pid: 1 is running!, global_variable: 9
[INTERRUPT] S mode timer interrupt!
[U-PARENT] pid: 1 is running!, global_variable: 10
[U-PARENT] pid: 1 is running!, global_variable: 11
[INTERRUPT] S mode timer interrupt!

switch to [PID = 2 COUNTER = 4 PRIORITY = 37]
[U-CHILD] pid: 2 is running!, global_variable: 3
[U-CHILD] pid: 2 is running!, global_variable: 4
[U-CHILD] pid: 2 is running!, global_variable: 5
[INTERRUPT] S mode timer interrupt!
[U-CHILD] pid: 2 is running!, global_variable: 6
[U-CHILD] pid: 2 is running!, global_variable: 7
[INTERRUPT] S mode timer interrupt!
[U-CHILD] pid: 2 is running!, global_variable: 8
[U-CHILD] pid: 2 is running!, global_variable: 9
[INTERRUPT] S mode timer interrupt!
[U-CHILD] pid: 2 is running!, global_variable: 10
[U-CHILD] pid: 2 is running!, global_variable: 11
[INTERRUPT] S mode timer interrupt!
SET [PID = 1 COUNTER = 1]
SET [PID = 2 COUNTER = 4]

switch to [PID = 1 COUNTER = 1 PRIORITY = 37]
[U-PARENT] pid: 1 is running!, global_variable: 12
[U-PARENT] pid: 1 is running!, global_variable: 13
[INTERRUPT] S mode timer interrupt!

```

第三个main函数：

```

2022 Hello RISC-V
ssstatus = 8000000000006000
sscratch = 0

switch to [PID = 1 COUNTER = 4 PRIORITY = 37]
Page Fault: sepc: 00000000000100e8, scause: 000000000000000c, stval: 00000000000100e8
Page Fault: sepc: 0000000000010158, scause: 000000000000000f, stval: 0000003fffffffff8
Page Fault: sepc: 0000000000010174, scause: 000000000000000d, stval: 0000000000011930
[U] pid: 1 is running!, global_variable: 0
[S] New task: 2
[U] pid: 1 is running!, global_variable: 1
[S] New task: 3
[U] pid: 1 is running!, global_variable: 2
[U] pid: 1 is running!, global_variable: 3
[U] pid: 1 is running!, global_variable: 4
[INTERRUPT] S mode timer interrupt!
[U] pid: 1 is running!, global_variable: 5
[U] pid: 1 is running!, global_variable: 6
[INTERRUPT] S mode timer interrupt!
[U] pid: 1 is running!, global_variable: 7
[U] pid: 1 is running!, global_variable: 8
[INTERRUPT] S mode timer interrupt!
[U] pid: 1 is running!, global_variable: 9
[U] pid: 1 is running!, global_variable: 10
[INTERRUPT] S mode timer interrupt!

switch to [PID = 2 COUNTER = 4 PRIORITY = 37]
[U] pid: 2 is running!, global_variable: 1
[S] New task: 4
[U] pid: 2 is running!, global_variable: 2
[U] pid: 2 is running!, global_variable: 3
[U] pid: 2 is running!, global_variable: 4
[INTERRUPT] S mode timer interrupt!
[U] pid: 2 is running!, global_variable: 5
[U] pid: 2 is running!, global_variable: 6
[INTERRUPT] S mode timer interrupt!
[U] pid: 2 is running!, global_variable: 7
[U] pid: 2 is running!, global_variable: 8
[INTERRUPT] S mode timer interrupt!
[U] pid: 2 is running!, global_variable: 9
[U] pid: 2 is running!, global_variable: 10
[INTERRUPT] S mode timer interrupt!

switch to [PID = 3 COUNTER = 4 PRIORITY = 37]
[U] pid: 3 is running!, global_variable: 2
[U] pid: 3 is running!, global_variable: 3
[U] pid: 3 is running!, global_variable: 4

```

## 更多测试样例

将INIT\_TASKS改为3（即初始化两个进程），测试结果如下：

```

switch to [PID = 2 COUNTER = 4 PRIORITY = 88]
[U-PARENT] pid: 2 is running! the 39th fibonacci number is 63245986 and the number @ 961 in the large array is 961

switch to [PID = 4 COUNTER = 4 PRIORITY = 88]
[U-CHILD] pid: 4 is running! the 39th fibonacci number is 63245986 and the number @ 961 in the large array is 961

switch to [PID = 3 COUNTER = 10 PRIORITY = 37]
[U-CHILD] pid: 3 is running! the 37th fibonacci number is 24157817 and the number @ 963 in the large array is 963
[U-CHILD] pid: 3 is running! the 38th fibonacci number is 39088169 and the number @ 962 in the large array is 962
[U-CHILD] pid: 3 is running! the 39th fibonacci number is 63245986 and the number @ 961 in the large array is 961
[U-CHILD] pid: 3 is running! the 40th fibonacci number is 102334155 and the number @ 960 in the large array is 960
SET [PID = 1 COUNTER = 10]
SET [PID = 2 COUNTER = 10]
SET [PID = 3 COUNTER = 5]
SET [PID = 4 COUNTER = 2]

switch to [PID = 4 COUNTER = 2 PRIORITY = 88]
[U-CHILD] pid: 4 is running! the 40th fibonacci number is 102334155 and the number @ 960 in the large array is 960

switch to [PID = 3 COUNTER = 5 PRIORITY = 37]

switch to [PID = 1 COUNTER = 10 PRIORITY = 37]
[U-PARENT] pid: 1 is running! the 37th fibonacci number is 24157817 and the number @ 963 in the large array is 963
[U-PARENT] pid: 1 is running! the 38th fibonacci number is 39088169 and the number @ 962 in the large array is 962
[U-PARENT] pid: 1 is running! the 39th fibonacci number is 63245986 and the number @ 961 in the large array is 961
[U-PARENT] pid: 1 is running! the 40th fibonacci number is 102334155 and the number @ 960 in the large array is 960

switch to [PID = 2 COUNTER = 10 PRIORITY = 88]
[U-PARENT] pid: 2 is running! the 40th fibonacci number is 102334155 and the number @ 960 in the large array is 960

```

## 思考题

1. 参考 `task_init` 创建一个新的 `task`，将的 `parent task` 的整个页复制到新创建的 `task_struct` 页上。这一步复制了哪些东西？

在这一步中，新创建的 `task_struct` 页复制了 `parent task` 的整个页。这包括了 `parent task` 的所有内存状态，如 `state` / `counter` / `priority` / `pid` / `thread` / `vma_cnt` / `vmass` 等等，以及 `parent task` 内核栈的所有信息（如 `pt_regs`），虽然之后我们会重新设置 `pid` 和 `thread` 中的内容，但是大部分状态是不变的。

2. 将 `thread.ra` 设置为 `__ret_from_fork`，并正确设置 `thread.sp`。仔细想想，这个应该设置成什么值？可以根据 `child task` 的返回路径来倒推。

`thread.sp` 应该设置为 `child task` 的 `pt_regs` 的地址，目前栈中有一个 `pt_regs`，`thread.sp` 指向 `child task` 的栈顶部，即现在的 `child_regs` 地址。

3. 利用参数 `regs` 来计算出 `child task` 的对应的 `pt_regs` 的地址，并将其中的 `a0`, `sp`, `sepc` 设置成正确的值。为什么还要设置 `sp`？

我之前的 `pt_regs` 中并没有保存 `sp`，所以不需要设置。若保存了 `sp`，则需要设置为当前内核栈，即 `child_regs` 地址