

练习0：填写已有实验

以下操作都是将内容复制到 lab5 里（都可以从 lab4 复制）**不要复制整个文件！！**

将 lab1 的 `kern/debug/kdebug.c`、`kern/init/init.c` 以及 `kern/trap/trap.c` 复制到 lab5 里

将 lab2 的 `kern/mm/pmm.c` 和 `kern/mm/default_pmm.c` 复制到 lab5 里

将 lab3 的 `kern/mm/vmm.c` 和 `kern/mm/swap_fifo.c` 复制到 lab5 里

最后将 lab4 的 `kern/process/proc.c` 复制到 lab5 里，之后还要改进代码

alloc_proc 函数源码

写于： `kern/process/proc.c`

```
1 // alloc_proc - alloc a proc_struct and init all fields of proc_struct
2 static struct proc_struct *
3 alloc_proc(void) {
4     struct proc_struct *proc = kmalloc(sizeof(struct proc_struct));
5     if (proc != NULL) {
6         //LAB4:EXERCISE1 YOUR CODE
7         /*
8          * below fields in proc_struct need to be initialized
9          *      enum proc_state state;           // Process state
10         *      int pid;                         // Process ID
11         *      int runs;                       // the running
12         times of Proces
13         *      uintptr_t kstack;               // Process kernel
14         stack
15         *      volatile bool need_resched;     // bool value:
16         need to be rescheduled to release CPU?
17         *      struct proc_struct *parent;     // the parent
18         process
19         *      struct mm_struct *mm;           // Process's
20         memory management field
21         *      struct context context;         // Switch here to
22         run process
23         *      struct trapframe *tf;          // Trap frame for
24         current interrupt
25         *      uintptr_t cr3;                 // CR3 register:
26         the base addr of Page Directroy Table(PDT)
27         *      uint32_t flags;                 // Process flag
28         *      char name[PROC_NAME_LEN + 1];  // Process name
29         */
30         //LAB5 YOUR CODE : (update LAB4 steps)
31         /*
32          * below fields(add in LAB5) in proc_struct need to be initialized
33          *      uint32_t wait_state;           // waiting state
34          *      struct proc_struct *cptr, *yptr, *optr; // relations
35         between processes
36         */
37     }
38     return proc;
39 }
```

可以看出来要多加两行初始化的代码

proc_struct 结构体

写于: kern/process/proc.h

先来看新的 proc_struct 结构体, 里面在最后多加了两行内容

```
1 struct proc_struct {
2     enum proc_state state;           // Process state
3     int pid;                         // Process ID
4     int runs;                        // the running times of
    Proces
5     uintptr_t kstack;               // Process kernel stack
6     volatile bool need_resched;     // bool value: need to be
    rescheduled to release CPU?
7     struct proc_struct *parent;     // the parent process
8     struct mm_struct *mm;           // Process's memory
    management field
9     struct context context;         // Switch here to run
    process
10    struct trapframe *tf;            // Trap frame for current
    interrupt
11    uintptr_t cr3;                   // CR3 register: the base
    addr of Page Directroy Table(PDT)
12    uint32_t flags;                  // Process flag
13    char name[PROC_NAME_LEN + 1];   // Process name
14    list_entry_t list_link;          // Process link list
15    list_entry_t hash_link;          // Process hash list
16    int exit_code;                   // exit code (be sent to
    parent proc)
17    uint32_t wait_state;              // waiting state
18    struct proc_struct *cptr, *yptr, *optr; // relations between
    processes
19 };
```

alloc_proc 函数答案

写于: kern/process/proc.h

```
1 static struct proc_struct *
2 alloc_proc(void) {
3     struct proc_struct *proc = kmalloc(sizeof(struct proc_struct));
4     if (proc != NULL) {
5         proc->pid = -1;               // 进程ID
6         memset(&(proc->name), 0, PROC_NAME_LEN); // 进程名
7         proc->state = PROC_UNINIT;    // 进程状态
8         proc->runs = 0;               // 进程时间片
9         proc->need_resched = 0;       // 进程是否能被
    调度
10        proc->flags = 0;               // 标志位
11        proc->kstack = 0;              // 进程所使用的
    内存栈地址
12        proc->cr3 = boot_cr3;         // 将页目录表地
    址设为内核页目录表基址
13        proc->mm = NULL;              // 进程所用的虚
    拟内存
    }
```

```

14     memset(&(proc->context), 0, sizeof(struct context)); // 进程的上下文
15     proc->tf = NULL; // 中断帧指针
16     proc->parent = NULL; // 该进程的父进
程
17     proc->wait_state = 0; // Lab5: 等待状
态的标志位
18     proc->cptr = NULL; // Lab5: 该进程
的子进程
19     proc->yptr = NULL; // Lab5: 该进程
的弟进程
20     proc->optr = NULL; // Lab5: 该进程
的兄进程
21 }
22 return proc;
23 }

```

do_fork 函数源码

写于: kern/process/proc.c

```

1  /* do_fork -      parent process for a new child process
2   * @clone_flags: used to guide how to clone the child process
3   * @stack:       the parent's user stack pointer. if stack==0, It means to
fork a kernel thread.
4   * @tf:          the trapframe info, which will be copied to child
process's proc->tf
5   */
6  int
7  do_fork(uint32_t clone_flags, uintptr_t stack, struct trapframe *tf) {
8      int ret = -E_NO_FREE_PROC;
9      struct proc_struct *proc;
10     if (nr_process >= MAX_PROCESS) {
11         goto fork_out;
12     }
13     ret = -E_NO_MEM;
14     //LAB4:EXERCISE2 YOUR CODE
15     /*
16      * Some Useful MACROS, Functions and DEFINES, you can use them in below
implementation.
17      * MACROS or Functions:
18      *   alloc_proc:   create a proc struct and init fields
(lab4:exercise1)
19      *   setup_kstack: alloc pages with size KSTACKPAGE as process kernel
stack
20      *   copy_mm:      process "proc" duplicate OR share process
"current"'s mm according clone_flags
21      *                 if clone_flags & CLONE_VM, then "share" ; else
"duplicate"
22      *   copy_thread:  setup the trapframe on the process's kernel stack
top and
23      *                 setup the kernel entry point and stack of process
24      *   hash_proc:    add proc into proc hash_list
25      *   get_pid:      alloc a unique pid for process
26      *   wakeup_proc:  set proc->state = PROC_RUNNABLE
27      * VARIABLES:
28      *   proc_list:    the process set's list
29      *   nr_process:   the number of process set

```

```

30     */
31
32     // 1. call alloc_proc to allocate a proc_struct
33     // 2. call setup_kstack to allocate a kernel stack for child process
34     // 3. call copy_mm to dup OR share mm according clone_flag
35     // 4. call copy_thread to setup tf & context in proc_struct
36     // 5. insert proc_struct into hash_list && proc_list
37     // 6. call wakeup_proc to make the new child process RUNNABLE
38     // 7. set ret vaule using child proc's pid
39
40     //LAB5 YOUR CODE : (update LAB4 steps)
41     /* Some Functions
42     *      set_links: set the relation links of process. ALSO SEE:
43     remove_links: lean the relation links of process
44     *      -----
45     *      update step 1: set child proc's parent to current process, make
46     sure current process's wait_state is 0
47     *      update step 5: insert proc_struct into hash_list && proc_list, set
48     the relation links of process
49     */
50
51 fork_out:
52     return ret;
53
54 bad_fork_cleanup_kstack:
55     put_kstack(proc);
56 bad_fork_cleanup_proc:
57     kfree(proc);
58     goto fork_out;
59 }

```

set_links 函数

写于: kern/process/proc.c

```

1 // set_links - set the relation links of process
2 static void
3 set_links(struct proc_struct *proc) {
4     list_add(&proc_list, &(proc->list_link)); // 将 proc 进程添加到进
        程列表里
5     proc->yptr = NULL; // 初始化 proc 进程的弟
        进程
6     if ((proc->optr = proc->parent->cptr) != NULL) { // 如果 proc 进程的兄进
        程等于 proc 进程的父进程的子进程
7         proc->optr->yptr = proc; // 就设置 proc 进程的兄
        进程的弟进程为 proc 进程
8     }
9     proc->parent->cptr = proc; // 使 proc 进程的父进程
        的弟进程为 proc 进程
10    nr_process ++; // 进程控制块总数加一
11 }

```

所以这个进程的主要目的就是 将 proc 进程控制块放到进程列表里

do_fork 函数答案

写于: kern/process/proc.c

因为 `set_links` 函数 里有 `nr_process ++` , 所以 `do_fork` 函数 里的 `nr_process ++` 就可以删掉了

```
1  int
2  do_fork(uint32_t clone_flags, uintptr_t stack, struct trapframe *tf) {
3      int ret = -E_NO_FREE_PROC;           // 设置返回值为 -5
4      struct proc_struct *proc;           // 定义 proc_struct 结构体
      指针变量 proc (子进程)
5      if (nr_process >= MAX_PROCESS) {    // 进程总数 (全局变量) 如果
      >= 0x1000
6          goto fork_out;                 // 就跳转到 fork_out 地址,
      此时返回值是 -5
7      }
8      ret = -E_NO_MEM;                   // 更换返回值为 -4
9      if ((proc = alloc_proc()) == NULL) { // 申请一个 proc_struct 结
      构体, 作为子进程的进程控制块
10         goto fork_out;                 // 若是申请失败就跳转到
      fork_out 地址, 此时返回值是 -4
11     }
12
13     proc->parent = current;             // 设置当前进程为上面新申请的
      子进程的父进程
14     assert(current->wait_state == 0);    // Lab5: 设置当前进程为等待
      进程
15
16     if (setup_kstack(proc) != 0) {      // 分配并初始化子进程的内核
      栈, 返回值是 0 则成功分配
17         goto bad_fork_cleanup_proc;     // 创建失败, 即内存不足, 则跳
      转到 bad_fork_cleanup_proc
18     }
19     if (copy_mm(clone_flags, proc) != 0) { // 将当前进程的内存信息复制到
      子进程
20         goto bad_fork_cleanup_kstack;   // 出错就跳转到
      bad_fork_cleanup_kstack
21     }
22     copy_thread(proc, stack, tf);        // 复制当前进程的中断帧和上下
      文到子进程中
23
24     bool intr_flag;
25     local_intr_save(intr_flag);         // 禁止中断发生, 保护代码运行
26     {
27         proc->pid = get_pid();           // 获取子进程的 pid
28         hash_proc(proc);                // 将子进程添加进哈希列表
29         set_links(proc);                // Lab5: 将子进程的进程控制
      块放到进程列表里
30     }
31     local_intr_restore(intr_flag);       // 如果 intr_flag 不为 0,
      则允许中断发生
32
33     wakeup_proc(proc);                  // 唤醒该进程
34
35     ret = proc->pid;                     // 更新返回值为子进程的 pid
36 fork_out:
37     return ret;
38
39 bad_fork_cleanup_kstack:
40     put_kstack(proc);                   // 释放子进程的内核栈
```

```

41 | bad_fork_cleanup_proc:
42 |     kfree(proc);                                // 释放申请的物理页，跳转到
    |     fork_out
43 |     goto fork_out;
44 | }

```

idt_init 函数源码

写于: kern/trap/trap.c

```

1 | /* idt_init - initialize IDT to each of the entry points in
   | kern/trap/vectors.S */
2 | void
3 | idt_init(void) {
4 |     /* LAB1 YOUR CODE : STEP 2 */
5 |     /* (1) where are the entry addrs of each Interrupt Service Routine
   | (ISR)?
6 |      *      All ISR's entry addrs are stored in __vectors. where is
   | uintptr_t __vectors[] ?
7 |      *      __vectors[] is in kern/trap/vector.S which is produced by
   | tools/vector.c
8 |      *      (try "make" command in lab1, then you will find vector.S in
   | kern/trap DIR)
9 |      *      You can use "extern uintptr_t __vectors[];" to define this
   | extern variable which will be used later.
10 |     /* (2) Now you should setup the entries of ISR in Interrupt
   | Description Table (IDT).
11 |     /*      Can you see idt[256] in this file? Yes, it's IDT! you can use
   | SETGATE macro to setup each item of IDT
12 |     /* (3) After setup the contents of IDT, you will let CPU know where is
   | the IDT by using 'lidt' instruction.
13 |     /*      You don't know the meaning of this instruction? just google it!
   | and check the libs/x86.h to know more.
14 |     /*      Notice: the argument of lidt is idt_pd. try to find it!
15 |     */
16 |     /* LAB5 YOUR CODE */
17 |     //you should update your lab1 code (just add ONE or TWO lines of
   | code), let user app to use syscall to get the service of ucore
18 |     //so you should setup the syscall interrupt gate in here
19 | }

```

idt_init 函数答案

写于: kern/trap/trap.c

```

1 void
2 idt_init(void) {
3     extern uintptr_t __vectors[];
4     int i;
5     for (i = 0; i < sizeof(idt) / sizeof(struct gatedesc); i++) {
6         SETGATE(idt[i], 0, GD_KTEXT, __vectors[i], DPL_KERNEL);
7     }
8     SETGATE(idt[T_SYSCALL], 1, GD_KTEXT, __vectors[T_SYSCALL], DPL_USER);
9     // Lab5
10    lidt(&idt_pd);
11 }

```

trap_dispatch 函数源码

写于: kern/trap/trap.c

```

1 static void
2 trap_dispatch(struct trapframe *tf) {
3     char c;
4     int ret = 0;
5
6     switch (tf->tf_trapno) {
7     case T_PGFLT: // page fault
8         if ((ret = pgfault_handler(tf)) != 0) {
9             print_trapframe(tf);
10            if (current == NULL) {
11                panic("handle pgfault failed. ret=%d\n", ret);
12            }
13            else {
14                if (trap_in_kernel(tf)) {
15                    panic("handle pgfault failed in kernel mode. ret=%d\n",
16ret);
17                }
18                cprintf("killed by kernel.\n");
19                panic("handle user mode pgfault failed. ret=%d\n", ret);
20                do_exit(-E_KILLED);
21            }
22        }
23        break;
24     case T_SYSCALL:
25         syscall();
26         break;
27     case IRQ_OFFSET + IRQ_TIMER:
28         #if 0
29             LAB3 : If some page replacement algorithm(such as CLOCK PRA) need tick
30             to change the priority of pages,
31             then you can add code here.
32         #endif
33         /* LAB1 YOUR CODE : STEP 3 */
34         /* handle the timer interrupt */
35         /* (1) After a timer interrupt, you should record this event using
36            a global variable (increase it), such as ticks in kern/driver/clock.c
37            * (2) Every TICK_NUM cycle, you can print some info using a
38            funciton, such as print_ticks().
39            * (3) Too Simple? Yes, I think so!
40         */

```

```

37      /* LAB5 YOUR CODE */
38      /* you should upate you lab1 code (just add ONE or TWO lines of
code):
39          *    Every TICK_NUM cycle, you should set current process's
current->need_resched = 1
40          */
41
42      break;
43      case IRQ_OFFSET + IRQ_COM1:
44          c = cons_getc();
45          cprintf("serial [%03d] %c\n", c, c);
46          break;
47      case IRQ_OFFSET + IRQ_KBD:
48          c = cons_getc();
49          cprintf("kbd [%03d] %c\n", c, c);
50          break;
51      // LAB1 CHALLENGE 1 : YOUR CODE you should modify below codes.
52      case T_SWITCH_TOU:
53      case T_SWITCH_TOK:
54          panic("T_SWITCH_** ??\n");
55          break;
56      case IRQ_OFFSET + IRQ_IDE1:
57      case IRQ_OFFSET + IRQ_IDE2:
58          /* do nothing */
59          break;
60      default:
61          print_trapframe(tf);
62          if (current != NULL) {
63              cprintf("unhandled trap.\n");
64              do_exit(-E_KILLED);
65          }
66          // in kernel, it must be a mistake
67          panic("unexpected trap in kernel.\n");
68
69      }
70  }

```

trap_dispatch 函数答案

写于: kern/trap/trap.c

```

1  static void
2  trap_dispatch(struct trapframe *tf) {
3      char c;
4      int ret = 0;
5
6      switch (tf->tf_trapno) {
7      case T_PGFLT: // page fault
8          if ((ret = pgfault_handler(tf)) != 0) {
9              print_trapframe(tf);
10             panic("handle pgfault failed. %e\n", ret);
11         }
12         break;
13      case IRQ_OFFSET + IRQ_TIMER:
14          ticks ++;
15          if (ticks % TICK_NUM == 0) { // 当时间片用完
16              assert(current != NULL); // Lab5: 判断 current 变量是否有值

```



```

17         current->need_resched = 1; // Lab5: 设置当前进程需要被调度
18     }
19     break;
20 case IRQ_OFFSET + IRQ_COM1:
21     c = cons_getc();
22     cprintf("serial [%03d] %c\n", c, c);
23     break;
24 case IRQ_OFFSET + IRQ_KBD:
25     c = cons_getc();
26     cprintf("kbd [%03d] %c\n", c, c);
27     break;
28 // LAB1 CHALLENGE 1 : YOUR CODE you should modify below codes.
29 case T_SWITCH_TOU:
30 case T_SWITCH_TOK:
31     panic("T_SWITCH_** ??\n");
32     break;
33 case IRQ_OFFSET + IRQ_IDE1:
34 case IRQ_OFFSET + IRQ_IDE2:
35     /* do nothing */
36     break;
37 default:
38     // in kernel, it must be a mistake
39     if ((tf->tf_cs & 3) == 0) {
40         print_trapframe(tf);
41         panic("unexpected trap in kernel.\n");
42     }
43 }
44 }

```

练习1: 加载应用程序并执行

mm_struct 结构体

```

1 // the control struct for a set of vma using the same PDT
2 struct mm_struct { // 描述一个进程的虚拟地址空间，每个进程的 pcb
    // 中会有一个指针指向本结构体
3     list_entry_t mmap_list; // 链接同一页目录表的虚拟内存空间中双向链表的头
    // 节点
4     struct vma_struct *mmap_cache; // 当前正在使用的虚拟内存空间
5     pde_t *pgdir; // mm_struct 所维护的页表地址(用来找 PTE)
6     int map_count; // 虚拟内存块数目
7     void *sm_priv; // 记录访问情况链表头地址(用于置换算法)
8     int mm_count; // 共享 mm 的进程数，或者说对使用该结构体的某一
    // 结构体变量的引用次数
9     lock_t mm_lock; // 互斥锁，用于使用 dup_mmap 复制 mm
10 };

```

这个结构体新加了 `mm_count` 和 `mm_lock` 变量

load_icode 函数源码

写于: `kern/process/proc.c`

```

1 /* load_icode - load the content of binary program(ELF format) as the new
   content of current process

```

```

2  * @binary: the memory addr of the content of binary program
3  * @size: the size of the content of binary program
4  */
5  static int
6  load_icode(unsigned char *binary, size_t size) {
7      if (current->mm != NULL) {
8          panic("load_icode: current->mm must be empty.\n");
9      }
10
11     int ret = -E_NO_MEM;
12     struct mm_struct *mm;
13     //(1) create a new mm for current process
14     if ((mm = mm_create()) == NULL) {
15         goto bad_mm;
16     }
17     //(2) create a new PDT, and mm->pgdir= kernel virtual addr of PDT
18     if (setup_pgdir(mm) != 0) {
19         goto bad_pgdir_cleanup_mm;
20     }
21     //(3) copy TEXT/DATA section, build BSS parts in binary to memory
    space of process
22     struct Page *page;
23     //(3.1) get the file header of the bianry program (ELF format)
24     struct elfhdr *elf = (struct elfhdr *)binary;
25     //(3.2) get the entry of the program section headers of the bianry
    program (ELF format)
26     struct proghdr *ph = (struct proghdr *) (binary + elf->e_phoff);
27     //(3.3) This program is valid?
28     if (elf->e_magic != ELF_MAGIC) {
29         ret = -E_INVALID_ELF;
30         goto bad_elf_cleanup_pgdir;
31     }
32
33     uint32_t vm_flags, perm;
34     struct proghdr *ph_end = ph + elf->e_phnum;
35     for (; ph < ph_end; ph++) {
36         //(3.4) find every program section headers
37         if (ph->p_type != ELF_PT_LOAD) {
38             continue ;
39         }
40         if (ph->p_filesz > ph->p_memsz) {
41             ret = -E_INVALID_ELF;
42             goto bad_cleanup_mmap;
43         }
44         if (ph->p_filesz == 0) {
45             continue ;
46         }
47         //(3.5) call mm_map fun to setup the new vma ( ph->p_va, ph->p_memsz)
48         vm_flags = 0, perm = PTE_U;
49         if (ph->p_flags & ELF_PF_X) vm_flags |= VM_EXEC;
50         if (ph->p_flags & ELF_PF_W) vm_flags |= VM_WRITE;
51         if (ph->p_flags & ELF_PF_R) vm_flags |= VM_READ;
52         if (vm_flags & VM_WRITE) perm |= PTE_W;
53         if ((ret = mm_map(mm, ph->p_va, ph->p_memsz, vm_flags, NULL)) !=
0) {
54             goto bad_cleanup_mmap;
55         }
56         unsigned char *from = binary + ph->p_offset;

```

```

57     size_t off, size;
58     uintptr_t start = ph->p_va, end, la = ROUNDDOWN(start, PGSIZE);
59
60     ret = -E_NO_MEM;
61
62     //(3.6) alloc memory, and copy the contents of every program section
63     (from, from+end) to process's memory (la, la+end)
64     end = ph->p_va + ph->p_filesz;
65     //(3.6.1) copy TEXT/DATA section of binary program
66     while (start < end) {
67         if ((page = pgdir_alloc_page(mm->pgdir, la, perm)) == NULL) {
68             goto bad_cleanup_mmap;
69         }
70         off = start - la, size = PGSIZE - off, la += PGSIZE;
71         if (end < la) {
72             size -= la - end;
73         }
74         memcpy(page2kva(page) + off, from, size);
75         start += size, from += size;
76     }
77
78     //(3.6.2) build BSS section of binary program
79     end = ph->p_va + ph->p_memsz;
80     if (start < la) {
81         /* ph->p_memsz == ph->p_filesz */
82         if (start == end) {
83             continue ;
84         }
85         off = start + PGSIZE - la, size = PGSIZE - off;
86         if (end < la) {
87             size -= la - end;
88         }
89         memset(page2kva(page) + off, 0, size);
90         start += size;
91         assert((end < la && start == end) || (end >= la && start ==
92             la));
93     }
94     while (start < end) {
95         if ((page = pgdir_alloc_page(mm->pgdir, la, perm)) == NULL) {
96             goto bad_cleanup_mmap;
97         }
98         off = start - la, size = PGSIZE - off, la += PGSIZE;
99         if (end < la) {
100             size -= la - end;
101         }
102         memset(page2kva(page) + off, 0, size);
103         start += size;
104     }
105
106     //(4) build user stack memory
107     vm_flags = VM_READ | VM_WRITE | VM_STACK;
108     if ((ret = mm_map(mm, USTACKTOP - USTACKSIZE, USTACKSIZE, vm_flags,
109         NULL)) != 0) {
110         goto bad_cleanup_mmap;
111     }
112     assert(pgdir_alloc_page(mm->pgdir, USTACKTOP-PGSIZE , PTE_USER) !=
113         NULL);

```

```

110     assert(pgdir_alloc_page(mm->pgdir, USTACKTOP-2*PGSIZE , PTE_USER) !=
        NULL);
111     assert(pgdir_alloc_page(mm->pgdir, USTACKTOP-3*PGSIZE , PTE_USER) !=
        NULL);
112     assert(pgdir_alloc_page(mm->pgdir, USTACKTOP-4*PGSIZE , PTE_USER) !=
        NULL);
113
114     //(5) set current process's mm, sr3, and set CR3 reg = physical addr
of Page Directory
115     mm_count_inc(mm);
116     current->mm = mm;
117     current->cr3 = PADDR(mm->pgdir);
118     lcr3(PADDR(mm->pgdir));
119
120     //(6) setup trapframe for user environment
121     struct trapframe *tf = current->tf;
122     memset(tf, 0, sizeof(struct trapframe));
123     /* LAB5:EXERCISE1 YOUR CODE
124      * should set tf_cs,tf_ds,tf_es,tf_ss,tf_esp,tf_eip,tf_eflags
125      * NOTICE: If we set trapframe correctly, then the user level process
can return to USER MODE from kernel. So
126      *         tf_cs should be USER_CS segment (see memlayout.h)
127      *         tf_ds=tf_es=tf_ss should be USER_DS segment
128      *         tf_esp should be the top addr of user stack (USTACKTOP)
129      *         tf_eip should be the entry point of this binary program
(elf->e_entry)
130      *         tf_eflags should be set to enable computer to produce
Interrupt
131      */
132     ret = 0;
133 out:
134     return ret;
135 bad_cleanup_mmap:
136     exit_mmap(mm);
137 bad_elf_cleanup_pgdir:
138     put_pgdir(mm);
139 bad_pgdir_cleanup_mm:
140     mm_destroy(mm);
141 bad_mm:
142     goto out;
143 }

```

mm_create 函数

写于: kern/mm/vmm.c

```

1 // mm_create - alloc a mm_struct & initialize it.
2 struct mm_struct *
3 mm_create(void) {
4     // 申请一块内存地址空间
5     struct mm_struct *mm = kcalloc(sizeof(struct mm_struct));
6
7     if (mm != NULL) {                                // 如果申请成功, 返回了内存地址空间
的地址
8         list_init(&(mm->mmap_list));                // 将该内存地址空间中的链表节点设置
为头节点
9         mm->mmap_cache = NULL;                      // 初始化虚拟内存空间

```

```

10      mm->pgdir = NULL;                                // 初始化 mm_struct 所维护的页表地
    址
11      mm->map_count = 0;                                // 初始化虚拟内存块的数目
12
13      if (swap_init_ok) swap_init_mm(mm); // 如果可以换入页面，那么就初始化用
    于置换算法的链表
14      else mm->sm_priv = NULL;                        // 否则就将记录访问情况的链表头地址
    设置为 NULL
15
16      set_mm_count(mm, 0);                             // 设置该虚拟内存块的数目为 0
17      lock_init(&(mm->mm_lock));                       // 初始化互斥锁
18  }
19  return mm;                                            // 返回 mm 内存地址空间的地址
20 }

```

该函数的作用是申请并初始化一块内存地址空间

swap_init_mm 函数

写于: kern/mm/swap.c

```

1  int
2  swap_init_mm(struct mm_struct *mm)
3  {
4      return sm->init_mm(mm);
5  }

```

- swap_manager_fifo 结构体变量

写于: kern/mm/swap_fifo.c

```

1  struct swap_manager swap_manager_fifo =
2  {
3      .name          = "fifo swap manager",
4      .init          = &_amp;_fifo_init,
5      .init_mm       = &_amp;_fifo_init_mm,
6      .tick_event    = &_amp;_fifo_tick_event,
7      .map_swappable = &_amp;_fifo_map_swappable,
8      .set_unswappable = &_amp;_fifo_set_unswappable,
9      .swap_out_victim = &_amp;_fifo_swap_out_victim,
10     .check_swap     = &_amp;_fifo_check_swap,
11 };

```

所以 init_mm 函数就是 _fifo_init_mm 函数

_fifo_init_mm 函数

写于: kern/mm/swap_fifo.c

```

1  /*
2   * (2) _fifo_init_mm: init pra_list_head and let mm->sm_priv point to the
   addr of pra_list_head.
3   *           Now, From the memory control struct mm_struct, we can
   access FIFO PRA
4   */
5   static int
6   _fifo_init_mm(struct mm_struct *mm)
7   {
8       list_init(&pra_list_head);
9       mm->sm_priv = &pra_list_head;
10      //cprintf(" mm->sm_priv %x in fifo_init_mm\n",mm->sm_priv);
11      return 0;
12  }

```

将 pra_list_head 的地址作为链表的头地址

将记录访问情况链表的头地址设置为 pra_list_head 的地址

set_mm_count 函数

写于: kern/mm/swap_fifo.c

```

1  static inline void
2  set_mm_count(struct mm_struct *mm, int val) {
3      mm->mm_count = val;
4  }

```

设置共享 mm 的进程数为 val

lock_init 函数

写于: kern/sync/sync.h

```

1  static inline void
2  lock_init(lock_t *lock) {
3      *lock = 0;
4  }

```

初始化互斥锁，将其设置为 0

setup_pgdir 函数

写于: kern/proc/proc.c

```

1  // setup_pgdir - alloc one page as PDT
2  static int
3  setup_pgdir(struct mm_struct *mm) {
4      struct Page *page;
5      if ((page = alloc_page()) == NULL) { // 申请一个物理页
6          return -E_NO_MEM;              // 申请失败就返回 -4
7      }
8      pde_t *pgdir = page2kva(page);      // 通过该物理页获取其内核虚拟地址（页目
   录表）
9      memcpy(pgdir, boot_pgdir, PGSIZE); // 将启动时页目录的虚拟地址中的内容复制到 pgdir 中

```

```

10 // 将 pgdir 中对应关于 VPT 的元素的内容设置为 pgdir 的地址转为物理地址后再加上标
    志位的值
11 pgdir[PDX(VPT)] = PADDR(pgdir) | PTE_P | PTE_W;
12 mm->pgdir = pgdir; // 设置 mm_struct 所维护的页表地址为
    pgdir 的地址（虚拟地址）
13 return 0;
14 }

```

该函数的作用就是创建一个页目录表，并准备好相应的初始化工作

VPT 宏

写于：kern/mm/memlayout.c

```

1 /* *
2  * virtual page table. Entry PDX[VPT] in the PD (Page Directory) contains
3  * a pointer to the page directory itself, thereby turning the PD into a page
4  * table, which maps all the PTEs (Page Table Entry) containing the page
    mappings
5  * for the entire virtual address space into that 4 Meg region starting at
    VPT.
6  * */
7 #define VPT                                0xFAC00000

```

虚拟页表的地址

ELF_MAGIC 宏

写于：libs/elf.h

```

1 #define ELF_MAGIC    0x464C457FU // "\x7FELF" in little endian

```

E_INVAL_ELF 宏

写于：libs/error.h

```

1 #define E_INVAL_ELF    8 // Invalid elf file

```

ELF_PT_LOAD 宏

写于：libs/elf.h

```

1 /* values for Proghdr::p_type */
2 #define ELF_PT_LOAD    1

```

表明该段可被加载

ELF_PF_X 宏

写于：libs/elf.h

```

1 /* flag bits for Proghdr::p_flags */
2 #define ELF_PF_X    1

```

表明有执行权限

ELF_PF_W 宏

写于: `libs/elf.h`

```
1 | #define ELF_PF_W 2
```

ELF_PF_R 宏

写于: `libs/elf.h`

```
1 | #define ELF_PF_R 4
```

mm_map 函数

写于: `kern/mm/vmm.c`

```
1  int
2  mm_map(struct mm_struct *mm, uintptr_t addr, size_t len, uint32_t vm_flags,
3         struct vma_struct **vma_store) {
4      // 规定页的起始地址和结束地址
5      uintptr_t start = ROUNDDOWN(addr, PGSIZE), end = ROUNDUP(addr + len,
6      PGSIZE);
7      if (!USER_ACCESS(start, end)) { // 判断当前地址是否能分配给用户使用
8          return -E_INVAL;           // 返回 -3
9      }
10     assert(mm != NULL);             // 要求传入的 mm 必须有地址，而不是申请失败后
    的空值
11
12     int ret = -E_INVAL;              // 设置返回值为 -3
13
14     struct vma_struct *vma;
15     // 更新 mm 里的 vma 并取得它的值，判断页的结束地址是否大于 vma 的起始地址
16     if ((vma = find_vma(mm, start)) != NULL && end > vma->vm_start) {
17         goto out;                   // 大于就有问题了，直接跳到 out
18     }
19     ret = -E_NO_MEM;                // 更新返回值为 -4
20
21     // 初始化 vma，其实经过上面的判断 vma 已经不可能是 NULL 了
22     if ((vma = vma_create(start, end, vm_flags)) == NULL) {
23         goto out;                   // 跳转到 out
24     }
25     insert_vma_struct(mm, vma);     // 将该 vma 对应的链表插入到 vma-
    >list_link 链表中
26     if (vma_store != NULL) {        // vma_store 如果不为空
27         *vma_store = vma;          // 更新 *vma_store 为 vma
28     }
29     ret = 0;                        // 更新返回值为 0
30
31 out:
32     return ret;
33 }
```

该函数的大致意思就是找到一个合法的 vma 内存，并更新 `mm->mmap_cache` 为这个新 vma

ROUNDUP 宏

写于: **libs/defs.h**

```
1  /* Round up to the nearest multiple of n */
2  #define ROUNDUP(a, n) ({
3      size_t __n = (size_t)(n);
4      (typeof(a))(ROUNDDOWN((size_t)(a) + __n - 1, __n));
5  })
```

注释意思：四舍五入操作，四舍五入到 $n + 1$ 的倍数

其实应该只要 n 不是 0, 都可以进行对于 a 的倍数的四舍五入

只是在这个 ucore 的代码里用的都是 2 的倍数 (都用的 PGSIZE == 4096)

如果 a 是 15 的话, `ROUNDUP(15, 4) == 16`

USER_ACCESS 宏

写于: kern/mm/memlayout.h

```
1 #define USER_ACCESS(start, end) \
2 (USERBASE <= (start) && (start) < (end) && (end) <= USERTOP)
```

判断当前地址是否能分配给用户使用

USERBASE 宏

写于: kern/mm/memlayout.h

```
1 #define USERBASE 0x00200000
```

用户地址的起始地址

USERTOP 宏

写于: kern/mm/memlayout.h

```
1 | #define USERTOP          0xB0000000
```

用户地址的结束地址

find_vma 函数

写于: kern/mm/vmm.c

```

1 // find_vma - find a vma (vma->vm_start <= addr <= vma->vm_end)
2 struct vma_struct *
3 find_vma(struct mm_struct *mm, uintptr_t addr) {
4     struct vma_struct *vma = NULL;
5     if (mm != NULL) {
6         vma = mm->mmap_cache; // 设置 vma 为该内存管
7                                // 理空间的虚拟内存空间
8         // 如果 vma 不为空且 addr 属于合法的地址
9         if (!(vma != NULL && vma->vm_start <= addr && vma->vm_end > addr))
10             {
11                 bool found = 0; // 初始化标记位为 0
12             }
13     }
14 }

```

```

10      list_entry_t *list = &(mm->mmap_list), *le = list;    // 获
    取头节点的信息
11      while ((le = list_next(le)) != list) { // 遍历双向链表
12          vma = le2vma(le, list_link);        // 依靠列表元素得到对应
    的 vma_struct 结构体地址
13          if (vma->vm_start<=addr && addr < vma->vm_end) { // 如
    果 addr 为合法地址
14              found = 1;                        // 改变标记位为 1，确
    认找到了符合的地址
15              break;                            // 退出循环
16          }
17      }
18      if (!found) {                                // 标记位为 0，说明没
    找到合法地址
19          vma = NULL;                            // 因为没找到合法地址，
    所以 vma 为 NULL
20      }
21  }
22      if (vma != NULL) {                            // 如果 vma 找到了合法
    地址
23          mm->mmap_cache = vma;                    // 更新该内存管理空间的
    虚拟内存空间为 vma
24      }
25  }
26      return vma;                                // 返回 vma
27  }

```

貌似就是去虚拟内存空间里找一块合法的空间来替代原来使用的空间

vma_create 函数

写于: kern/mm/vmm.c

```

1  // vma_create - alloc a vma_struct & initialize it. (addr range:
    vm_start~vm_end)
2  struct vma_struct *
3  vma_create(uintptr_t vm_start, uintptr_t vm_end, uint32_t vm_flags) {
4      struct vma_struct *vma = kcalloc(sizeof(struct vma_struct));
5
6      if (vma != NULL) {
7          vma->vm_start = vm_start;
8          vma->vm_end = vm_end;
9          vma->vm_flags = vm_flags;
10     }
11     return vma;
12 }

```

用来初始化 vma 的函数

insert_vma_struct 函数

写于: kern/mm/vmm.c

```

1  // insert_vma_struct -insert vma in mm's list link
2  void
3  insert_vma_struct(struct mm_struct *mm, struct vma_struct *vma) {

```

```

4      assert(vma->vm_start < vma->vm_end);           // 规定起始地址必
    须小于结束地址
5      list_entry_t *list = &(mm->mmap_list);         // 获得虚拟内存空
    间的头节点
6      list_entry_t *le_prev = list, *le_next;        // 初始化前驱节点
    和后继节点
7
8      list_entry_t *le = list;                       // 初始化链表节点
    元素
9      while ((le = list_next(le)) != list) {         // 遍历双向链表
10         // 依靠列表元素得到对应的 vma_struct 结构体地址
11         struct vma_struct *mmap_prev = le2vma(le, list_link);
12         if (mmap_prev->vm_start > vma->vm_start) {   // 如果该新虚拟内
    存空间的起始地址大于原始的起始地址
13             break;                                   // 退出循环
14         }
15         le_prev = le;                               // 设置上一个列表
    元素为 le
16     }
17
18     le_next = list_next(le_prev);                   // le 的下一个元
    素为 le_next
19
20     /* check overlap */
21     if (le_prev != list) {                           // 如果 le 不等
    于 list
22         check_vma_overlap(le2vma(le_prev, list_link), vma); // 检查是否出现重
    叠的内存区域
23     }
24     if (le_next != list) {                           // 如果 le_next
    不等于 list
25         check_vma_overlap(vma, le2vma(le_next, list_link)); // 检查是否出现重
    叠的内存区域
26     }
27
28     vma->vm_mm = mm;                                 // 更新虚拟内存空
    间属于的内存管理区域
29     list_add_after(le_prev, &(vma->list_link));      // 将 le 元素添
    加到 vma->list_link 链表中
30
31     mm->map_count ++;                                // 虚拟内存块的数
    目自增 1
32 }

```

这个函数的作用就是将该 vma 对应的链表插入到 vma->list_link 链表中

check_vma_overlap 函数

写于: kern/mm/vmm.c

```

1 // check_vma_overlap - check if vma1 overlaps vma2 ?
2 static inline void
3 check_vma_overlap(struct vma_struct *prev, struct vma_struct *next) {
4     assert(prev->vm_start < prev->vm_end);
5     assert(prev->vm_end <= next->vm_start);
6     assert(next->vm_start < next->vm_end);
7 }

```

检查分配的内存区域是否出现重叠情况，出现就报错中止程序

正常情况下，起始地址都要小于相应的结束地址，且 prev 的结束地址要小于 next 的起始地址

load_icode 函数答案

写于: kern/process/proc.c

```
1 static int
2 load_icode(unsigned char *binary, size_t size) {
3     if (current->mm != NULL) { // 判断当前进程
4         panic("load_icode: current->mm must be empty.\n"); // 不为空就报错
5         // 退出程序
6     }
7     int ret = -E_NO_MEM; // 设置返回值为
8     struct mm_struct *mm; // 创建内存地址
9     // (1) create a new mm for current process
10    if ((mm = mm_create()) == NULL) { // 申请并初始化
11        // 一块内存地址空间
12        goto bad_mm; // 申请失败就跳
13    } // (2) create a new PDT, and mm->pgdir= kernel virtual addr of PDT
14    if (setup_pgdir(mm) != 0) { // 创建一个页目
15        // 录表
16        goto bad_pgdir_cleanup_mm; // 创建失败就跳
17    } // (3) copy TEXT/DATA section, build BSS parts in binary to memory
18    // space of process
19    struct Page *page;
20    // (3.1) get the file header of the binary program (ELF format)
21    struct elfhdr *elf = (struct elfhdr *)binary;
22    // (3.2) get the entry of the program section headers of the binary
23    // program (ELF format)
24    struct proghdr *ph = (struct proghdr *) (binary + elf->e_phoff);
25    // (3.3) This program is valid?
26    if (elf->e_magic != ELF_MAGIC) { // 检查程序文件
27        // 头是否合法
28        ret = -E_INVALID ELF; // 更改返回值为
29        // -8
30        goto bad_elf_cleanup_pgdir; // 跳转到
31    } bad_elf_cleanup_pgdir
32    }
33    uint32_t vm_flags, perm;
34    struct proghdr *ph_end = ph + elf->e_phnum;
35    for (; ph < ph_end; ph++) {
36        // (3.4) find every program section headers
37        if (ph->p_type != ELF_PT_LOAD) { // 如果该段不是
38            // 可加载的段
39            continue; // 就跳过
40        }
41    }
```

```

36         if (ph->p_filesz > ph->p_memsz) { // 如果段的大小
        大于段的内存大小
37             ret = -E_INVALID_ELF; // 更改返回值为
        -8
38             goto bad_cleanup_mmap; // 跳转到
        bad_cleanup_mmap
39         }
40         if (ph->p_filesz == 0) { // 如果段的大小
        等于 0
41             continue ; // 就跳过
42         }
43         //(3.5) call mm_map fun to setup the new vma ( ph->p_va, ph->p_memsz)
44         vm_flags = 0, perm = PTE_U;
45         if (ph->p_flags & ELF_PF_X) vm_flags |= VM_EXEC; // 如果文件有执
        行权限, 将虚拟内存空间页设置为有执行权限
46         if (ph->p_flags & ELF_PF_W) vm_flags |= VM_WRITE; // 如果文件有可
        写权限, 将虚拟内存空间页设置为有可写权限
47         if (ph->p_flags & ELF_PF_R) vm_flags |= VM_READ; // 如果文件有可
        读权限, 将虚拟内存空间页设置为有可读权限
48         if (vm_flags & VM_WRITE) perm |= PTE_W; // 如果虚拟内存
        空间页有可写权限, 将页目录表设置为有可写权限
49         // 找到一个合法的 vma 内存, 并更新 mm->mmap_cache 为这个新 vma
50         if ((ret = mm_map(mm, ph->p_va, ph->p_memsz, vm_flags, NULL)) !=
        0) {
51             goto bad_cleanup_mmap; // 没找到合法的
        vma 就跳转到 bad_cleanup_mmap
52         }
53         unsigned char *from = binary + ph->p_offset;
54         size_t off, size;
55         uintptr_t start = ph->p_va, end, la = ROUNDDOWN(start, PGSIZE);
56
57         ret = -E_NO_MEM; // 更新返回值为
        -4
58
59         //(3.6) alloc memory, and copy the contents of every program section
        (from, from+end) to process's memory (la, la+end)
60         end = ph->p_va + ph->p_filesz;
61         //(3.6.1) copy TEXT/DATA section of binary program
62         while (start < end) { // 如果 start
        小于 end 就一直遍历
63             if ((page = pgdir_alloc_page(mm->pgdir, la, perm)) == NULL) {
        // 申请一个页表
64                 goto bad_cleanup_mmap; // 跳转到
        bad_cleanup_mmap
65             }
66             off = start - la, size = PGSIZE - off, la += PGSIZE;
67             if (end < la) { // 如果 end 小
        于 la
68                 size -= la - end; // 用 size 减去
        (la + end)
69             }
70             memcpy(page2kva(page) + off, from, size); // 将 from 的内
        容复制到 page 对应的虚拟地址中
71             start += size, from += size; // 更新 start
        和 from 的地址
72         }
73
74         //(3.6.2) build BSS section of binary program

```

```

75     end = ph->p_va + ph->p_memsz;
76     if (start < 1a) {
77         /* ph->p_memsz == ph->p_filesz */
78         if (start == end) {
79             continue ;
80         }
81         off = start + PGSIZE - 1a, size = PGSIZE - off;
82         if (end < 1a) {
83             size -= 1a - end;
84         }
85         memset(page2kva(page) + off, 0, size);
86         start += size;
87         assert((end < 1a && start == end) || (end >= 1a && start ==
1a));
88     }
89     while (start < end) {
90         if ((page = pgdir_alloc_page(mm->pgdir, 1a, perm)) == NULL) {
91             goto bad_cleanup_mmap;
92         }
93         off = start - 1a, size = PGSIZE - off, 1a += PGSIZE;
94         if (end < 1a) {
95             size -= 1a - end;
96         }
97         memset(page2kva(page) + off, 0, size);
98         start += size;
99     }
100 }
101 // (4) build user stack memory
102 vm_flags = VM_READ | VM_WRITE | VM_STACK;
103 if ((ret = mm_map(mm, USTACKTOP - USTACKSIZE, USTACKSIZE, vm_flags,
NULL)) != 0) {
104     goto bad_cleanup_mmap;
105 }
106 assert(pgdir_alloc_page(mm->pgdir, USTACKTOP-PGSIZE , PTE_USER) !=
NULL);
107 assert(pgdir_alloc_page(mm->pgdir, USTACKTOP-2*PGSIZE , PTE_USER) !=
NULL);
108 assert(pgdir_alloc_page(mm->pgdir, USTACKTOP-3*PGSIZE , PTE_USER) !=
NULL);
109 assert(pgdir_alloc_page(mm->pgdir, USTACKTOP-4*PGSIZE , PTE_USER) !=
NULL);
110
111 // (5) set current process's mm, sr3, and set CR3 reg = physical addr
of Page Directory
112 mm_count_inc(mm);
113 current->mm = mm;
114 current->cr3 = PADDR(mm->pgdir);
115 lcr3(PADDR(mm->pgdir));
116
117 // (6) setup trapframe for user environment
118 struct trapframe *tf = current->tf;
119 memset(tf, 0, sizeof(struct trapframe));
120 // Lab5
121 tf->tf_cs = USER_CS;
122 tf->tf_ds = tf->tf_es = tf->tf_ss = USER_DS;
123 tf->tf_esp = USTACKTOP;
124 tf->tf_eip = elf->e_entry;
125 tf->tf_eflags = FL_IF;

```

```

126     ret = 0;
127 out:
128     return ret;
129 bad_cleanup_mmap:
130     exit_mmap(mm);
131 bad_elf_cleanup_pgdir:
132     put_pgdir(mm);
133 bad_pgdir_cleanup_mm:
134     mm_destroy(mm);
135 bad_mm:
136     goto out;
137 }

```

分析不下去了，反正这部分就差个布置用户态的栈信息的代码，加上就好了

练习2: 父进程复制自己的内存空间给子进程

copy_range 函数源码

写于: kern/mm/pmm.c

```

1  int
2  copy_range(pde_t* to, pde_t* from, uintptr_t start, uintptr_t end, bool
share){
3      assert(start % PGSIZE == 0 && end % PGSIZE == 0);
4      assert(USER_ACCESS(start, end));
5      // copy content by page unit.
6      do{
7          //call get_pte to find process A's pte according to the addr start
8          pte_t* ptep = get_pte(from, start, 0), * nptep;
9          if(ptep == NULL){
10             start = ROUNDDOWN(start + PTSIZE, PTSIZE);
11             continue;
12         }
13         //call get_pte to find process B's pte according to the addr start. If
pte is NULL, just alloc a PT
14         if(*ptep & PTE_P){
15             if((nptep = get_pte(to, start, 1)) == NULL){
16                 return -E_NO_MEM;
17             }
18             uint32_t perm = (*ptep & PTE_USER);
19             //get page from ptep
20             struct Page* page = pte2page(*ptep);
21             // alloc a page for process B
22             struct Page* npage = alloc_page();
23             assert(page != NULL);
24             assert(npage != NULL);
25             int ret = 0;
26             /* LAB5:EXERCISE2 YOUR CODE
27              * replicate content of page to npage, build the map of phy
addr of nage with the linear addr start
28              *
29              * Some Useful MACROS and DEFINES, you can use them in below
implementation.
30              * MACROS or Functions:

```

```

31         *    page2kva(struct Page *page): return the kernel virtual
addr of memory which page managed (SEE pmm.h)
32         *    page_insert: build the map of phy addr of an Page with
the linear addr la
33         *    memcpy: typical memory copy function
34         *
35         * (1) find src_kvaddr: the kernel virtual address of page
36         * (2) find dst_kvaddr: the kernel virtual address of npage
37         * (3) memory copy from src_kvaddr to dst_kvaddr, size is
PGSIZE
38         * (4) build the map of phy addr of nage with the linear addr
start
39         */
40         assert(ret == 0);
41     }
42     start += PGSIZE;
43 } while(start != 0 && start < end);
44 return 0;
45 }

```

copy_range 函数答案

写于: kern/mm/pmm.c

```

1  int
2  copy_range(pde_t* to, pde_t* from, uintptr_t start, uintptr_t end, bool
share){
3      assert(start % PGSIZE == 0 && end % PGSIZE == 0);    // 确保 start 和
end 的大小都是页对齐的
4      assert(USER_ACCESS(start, end));                      // 确保 start 和
end 是在用户态空间范围内的
5      // copy content by page unit.
6      do{
7          //call get_pte to find process A's pte according to the addr start
8          pte_t* ptep = get_pte(from, start, 0), * nptep;    // 获得一个有存在位
标志的 ptep
9          if(ptep == NULL){                                  // 如果没有成功获得
有存在位标志的 ptep
10             start = ROUNDDOWN(start + PTSIZE, PTSIZE);    // 将 start 设置为
start + PTSIZE 的大小
11             continue;                                       // 继续下一个循环
12         }
13         //call get_pte to find process B's pte according to the addr start. If
pte is NULL, just alloc a PT
14         if(*ptep & PTE_P){                                  // 如果得到的 ptep
有存在位标志
15             if((nptep = get_pte(to, start, 1)) == NULL){ // 新分配一个页表
16                 return -E_NO_MEM;                          // 分配失败就返回 -4
17             }
18             uint32_t perm = (*ptep & PTE_USER);            // 设置用户标志位(用
户位、可写位、存在位)
19             //get page from ptep
20             struct Page* page = pte2page(*ptep);          // 从页表得到对应的
物理页
21             // alloc a page for process B
22             struct Page* npage = alloc_page();            // 申请一块物理页

```



```

23         assert(page != NULL); // 没得到对应的物理
    页, 中止程序
24         assert(npag e != NULL); // 申请物理页不成
    功, 中止程序
25         int ret = 0; // 设置返回值为 0
26
27         void* kva_src = page2kva(page); // Lab5: 得到 page 的虚拟地址
28         void* kva_dst = page2kva(npag e); // Lab5: 得到 npag e 的虚拟地址
29         memcpy(kva_dst, kva_src, PGSIZE); // Lab5: 将 page 虚拟地址里的内
    容复制到 npag e 的虚拟地址里
30         ret = page_insert(to, npag e, start, perm); // Lab5: 释放原先的二级
    页表映射并建立和 npag e 映射的页表
31         assert(ret == 0); // 返回值为 0 代表函
    数运行正常, 返回 -4 就中止程序
32     }
33     start += PGSIZE; // 将起始地址加上一
    个页的大小
34     } while(start != 0 && start < end); // 如果 start 不大
    于等于 end 且不等于 0 就一直循环
35     return 0;
36 }

```

练习3: 阅读分析源代码，理解进程执行 fork/exec/wait/exit 的实现，以及系统调用的实现

do_fork 函数分析

写于: kern/process/proc.c

```

1  int
2  do_fork(uint32_t clone_flags, uintptr_t stack, struct trapframe *tf) {
3      int ret = -E_NO_FREE_PROC; // 设置返回值为 -5
4      struct proc_struct *proc; // 定义 proc_struct 结构体指针变量 proc (子进程)
5      if (nr_process >= MAX_PROCESS) { // 进程总数 (全局变量) 如果 >= 0x1000
6          goto fork_out; // 就跳转到 fork_out 地址, 此时返回值是 -5
7      }
8      ret = -E_NO_MEM; // 更换返回值为 -4
9      if ((proc = alloc_proc()) == NULL) { // 申请一个 proc_struct 结构体
10         goto fork_out; // 若是申请失败就跳转到 fork_out 地址, 此时返回值是 -4
11     }
12
13     proc->parent = current; // 设置当前进程的父进程地址
14
15     if (setup_kstack(proc) != 0) { // 分配并初始化内核栈, 返回值是 0 则成功分配
16         // 创建失败, 即内存不足, 则跳转到 bad_fork_cleanup_proc
17         goto bad_fork_cleanup_proc;
18     }
19     if (copy_mm(clone_flags, proc) != 0) { // 将父进程的内存信息复制到子进程
20         // 出错就跳转到 bad_fork_cleanup_kstack
21         goto bad_fork_cleanup_kstack;
22     }
23     copy_thread(proc, stack, tf); // 复制父进程的中断帧和上下文
24
25     bool intr_flag;
26     local_intr_save(intr_flag); // 禁止中断发生, 保护代码运行
27     {

```

```

28     proc->pid = get_pid();           // 获取该程序的 pid
29     hash_proc(proc);                // 将该节点添加进哈希列表
30     list_add(&proc_list, &(proc->list_link)); // 将该节点添加进双向链表
31     nr_process ++;                  // 记录总进程数的变量自增 1
32 }
33 local_intr_restore(intr_flag);      // 如果 intr_flag 不为 0, 则允许中断发生
34
35 wakeup_proc(proc);                  // 唤醒该进程
36
37 ret = proc->pid;                     // 更新返回值为该进程的 pid
38 fork_out:
39     return ret;
40
41 bad_fork_cleanup_kstack:
42     put_kstack(proc);                // 释放内核栈
43 bad_fork_cleanup_proc:
44     kfree(proc);                     // 释放申请的物理页, 跳转到 fork_out
45     goto fork_out;
46 }

```

1. 分配并初始化进程控制块 (alloc_proc 函数)
2. 分配并初始化内核栈 (setup_stack 函数)
3. 根据 clone_flag 标志复制或共享进程内存管理结构 (copy_mm 函数)
4. 设置进程在内核 (将来也包括用户态) 正常运行和调度所需的中断帧和执行上下文 (copy_thread 函数)
5. 把设置好的进程控制块放入 hash_list 和 proc_list 两个全局进程链表中
6. 自此, 进程已经准备好执行了, 把进程状态设置为“就绪”态
7. 设置返回码为子进程的 id 号

do_exit 函数分析

写于: kern/process/proc.c

```

1 // do_exit - called by sys_exit
2 // 1. call exit_mmap & put_pgdir & mm_destroy to free the almost all
   memory space of process
3 // 2. set process' state as PROC_ZOMBIE, then call wakeup_proc(parent) to
   ask parent reclaim itself.
4 // 3. call scheduler to switch to other process
5 int
6 do_exit(int error_code) {
7     if (current == idleproc) {
8         panic("idleproc exit.\n");
9     }
10    if (current == initproc) {
11        panic("initproc exit.\n");
12    }
13
14    struct mm_struct *mm = current->mm;
15    if (mm != NULL) {
16        lcr3(boot_cr3);
17        if (mm_count_dec(mm) == 0) {
18            exit_mmap(mm);
19            put_pgdir(mm);
20            mm_destroy(mm);
21        }

```

```

22     current->mm = NULL;
23 }
24 current->state = PROC_ZOMBIE;
25 current->exit_code = error_code;
26
27 bool intr_flag;
28 struct proc_struct *proc;
29 local_intr_save(intr_flag);
30 {
31     proc = current->parent;
32     if (proc->wait_state == WT_CHILD) {
33         wakeup_proc(proc);
34     }
35     while (current->cptr != NULL) {
36         proc = current->cptr;
37         current->cptr = proc->optr;
38
39         proc->yptr = NULL;
40         if ((proc->optr == initproc->cptr) != NULL) {
41             initproc->cptr->yptr = proc;
42         }
43         proc->parent = initproc;
44         initproc->cptr = proc;
45         if (proc->state == PROC_ZOMBIE) {
46             if (initproc->wait_state == WT_CHILD) {
47                 wakeup_proc(initproc);
48             }
49         }
50     }
51 }
52 local_intr_restore(intr_flag);
53
54 schedule();
55 panic("do_exit will not return!! %d.\n", current->pid);
56 }

```

1. 先判断是否是用户进程，如果是，则开始回收此用户进程所占用的用户态虚拟内存空间（具体的回收过程不作详细说明）
2. 设置当前进程的终止性状态为 PROC_ZOMBIE，然后设置当前进程的退出码为 error_code
表明此时这个进程已经无法再被调用了，只能等待父进程来完成最后的回收工作（主要是回收该子进程的内核栈、进程控制块）
3. 如果当前父进程已经处于等待子进程的状态，即父进程的 wait_state 被置为 WT_CHILD
则此时就可以唤醒父进程，让父进程来帮子进程完成最后的资源回收工作。
4. 如果当前进程还有子进程，则需要把这些子进程的父进程指针设置为内核线程 init，且各个子进程指针需要插入到 init 的子进程链表中。如果某个子进程的执行状态是 PROC_ZOMBIE，
则需要唤醒 init 来完成对此子进程的最后回收工作
5. 执行 schedule() 调度函数，选择新的进程执行

do_execve 函数分析

写于：kern/process/proc.c

```

1 // do_execve - call exit_mmap(mm)&put_pgdir(mm) to reclaim memory space of
  // current process
2 // - call load_icode to setup new memory space according binary
  // prog.
3 int
4 do_execve(const char *name, size_t len, unsigned char *binary, size_t size)
  {
5     struct mm_struct *mm = current->mm;
6     if (!user_mem_check(mm, (uintptr_t)name, len, 0)) {
7         return -E_INVALID;
8     }
9     if (len > PROC_NAME_LEN) {
10         len = PROC_NAME_LEN;
11     }
12
13     char local_name[PROC_NAME_LEN + 1];
14     memset(local_name, 0, sizeof(local_name));
15     memcpy(local_name, name, len);
16
17     if (mm != NULL) {
18         lcr3(boot_cr3);
19         if (mm_count_dec(mm) == 0) {
20             exit_mmap(mm);
21             put_pgdir(mm);
22             mm_destroy(mm);
23         }
24         current->mm = NULL;
25     }
26     int ret;
27     if ((ret = load_icode(binary, size)) != 0) {
28         goto execve_exit;
29     }
30     set_proc_name(current, local_name);
31     return 0;
32
33 execve_exit:
34     do_exit(ret);
35     panic("already exit: %e.\n", ret);
36 }

```

1. 首先为加载新的执行码做好用户态内存空间清空准备。如果 mm 不为 NULL，则设置页表为内核空间页表，
且进一步判断 mm 的引用计数减 1 后是否为 0；如果为 0，则表明没有进程再需要此进程所占用的内存空间，
为此将根据 mm 中的记录，释放进程所占用户空间内存和进程页表本身所占空间，最后把当前进程的 mm 内存管理指针为空
2. 接下来是加载应用程序执行码到当前进程的新创建的用户态虚拟空间中，之后就是调用 load_icode 从而使之准备好执行

do_wait 函数分析

写于: kern/process/proc.c

```

1 // do_wait - wait one OR any children with PROC_ZOMBIE state, and free
  // memory space of kernel stack

```

```

2 //      - proc struct of this child.
3 // NOTE: only after do_wait function, all resources of the child proces are
  free.
4 int
5 do_wait(int pid, int *code_store) {
6     struct mm_struct *mm = current->mm;
7     if (code_store != NULL) {
8         if (!user_mem_check(mm, (uintptr_t)code_store, sizeof(int), 1)) {
9             return -E_INVALID;
10        }
11    }
12
13    struct proc_struct *proc;
14    bool intr_flag, haskid;
15    repeat:
16        haskid = 0;
17        if (pid != 0) {
18            proc = find_proc(pid);
19            if (proc != NULL && proc->parent == current) {
20                haskid = 1;
21                if (proc->state == PROC_ZOMBIE) {
22                    goto found;
23                }
24            }
25        }
26        else {
27            proc = current->cptr;
28            for (; proc != NULL; proc = proc->optr) {
29                haskid = 1;
30                if (proc->state == PROC_ZOMBIE) {
31                    goto found;
32                }
33            }
34        }
35        if (haskid) {
36            current->state = PROC_SLEEPING;
37            current->wait_state = WT_CHILD;
38            schedule();
39            if (current->flags & PF_EXITING) {
40                do_exit(-E_KILLED);
41            }
42            goto repeat;
43        }
44        return -E_BAD_PROC;
45
46    found:
47        if (proc == idleproc || proc == initproc) {
48            panic("wait idleproc or initproc.\n");
49        }
50        if (code_store != NULL) {
51            *code_store = proc->exit_code;
52        }
53        local_intr_save(intr_flag);
54        {
55            unhash_proc(proc);
56            remove_links(proc);
57        }
58        local_intr_restore(intr_flag);

```

```
59     put_kstack(proc);
60     kfree(proc);
61     return 0;
62 }
```

1. 如果 `pid != 0`，表示只找一个进程 id 号为 `pid` 的退出状态的子进程，否则找任意一个处于退出状态的子进程
2. 如果此子进程的执行状态不为 `PROC_ZOMBIE`，表明此子进程还没有退出，则当前进程设置执行状态为 `PROC_SLEEPING`（睡眠）
睡眠原因为 `WT_CHILD`（即等待子进程退出），调用 `schedule()` 函数选择新的进程执行，自己睡眠等待
如果被唤醒，则重复跳回步骤 1 处执行
3. 如果此子进程的执行状态为 `PROC_ZOMBIE`，表明此子进程处于退出状态，
需要当前进程（即子进程的父进程）完成对子进程的最终回收工作，
即首先把子进程控制块从两个进程队列 `proc_list` 和 `hash_list` 中删除，并释放子进程的内核堆栈和进程控制块。
自此，子进程才彻底地结束了它的执行过程，它所占用的所有资源均已释放。

扩展练习 Challenge：实现 Copy on Write（COW）机制

咕咕咕