

Report for lab4, YiChen Mao

Report for lab4, YiChen Mao

Part A: Multiprocessor Support and Cooperative Multitasking

Exercise 1.

Exercise 2.

Exercise 3.

Exercise 4.

Exercise 5.

Exercise 6.

Exercise 7.

Part B: Copy-on-Write Fork

Exercise 8.

Exercise 9.

Exercise 10.

Exercise 11.

Exercise 12.

Part C: Preemptive Multitasking and Inter-Process communication (IPC)

Exercise 13.

Exercise 14.

Exercise 15.

Challenge

Part A: Multiprocessor Support and Cooperative Multitasking

Exercise 1.

Task. 实现 kern/pmap.c 中的函数 mmio_map_region。

- 1 主要任务是把物理地址 `[pa, pa+size)` 映射到虚拟地址 `[MMIOBASE, MMIOBASE+size)` 处, 这里 `page` 可能并非页对齐的, 因此首先将 `pa` 和 `page` 进行页对齐 (注意到映射 `[pa, pa+size)`, 因此是将 `pa+size` 上取对齐), 由于 `base` 每次更新都是页对齐的因此若 `base` 未对齐一定出现错误。
- 2 然后由于 IO 地址空间只有 `1<<20` 字节, 判断是否出界, 然后用 `boot_map_region` 建立映射, 根据提示修改 权限位。
- 3 最后修正 `base`, (如果不修改 `base` 那么每次分配的空间初始位置都相同无法通过 `check_page`)

```

4  并返回映射的起始虚拟地址。
5
6  void *
7  mmio_map_region(physaddr_t pa, size_t size)
8  {
9      static uintptr_t base = MMIOBASE;
10     if (PGOFF(base))
11         panic("mmio_map_region: base error!");
12     size = ROUNDUP(size, PGSIZE);
13     if (size > PTSIZE || base + size >= MMIOLIM)
14         panic("mmio_map_region: error!");
15     boot_map_region(kern_pgdir, base, size, pa, PTE_PCD | PTE_PWT | PTE_W);
16     base += size;
17     return (void*)(base-size);
18     panic("mmio_map_region not implemented");
19 }

```

Exercise 2.

Task. 阅读 kern/init.c 中的函数 boot_aps() 和 mp_main(), 阅读 kern/mpentry.S 中的汇编, 并修改 kern/pmap.c 中 page_init 的实现, 已达到不把 MPENTRY_PADDR 所对应段加入空闲链表中。

```

1  mp_main() 应该是获取一些硬件信息和初始化一些结构。(为后面操作做准备)
2  boot_aps() 是先将 mpentry.S 中汇编代码复制到 MPENTRY_PADDR(0x7000) 并启动其他
   处理器, 其中调用了 mpentry.S 中的代码。
3  mpentry.S 初始化页表, 建立栈等(环境已经在 mp_main 之间建好)
4
5  只需要特殊判读一下 MPENTRY_PADDR, 较为简单, 具体代码见 kern/pmap.c。

```

Question.

比较 kern/mpentry.S 和 boot/boot.S, 为什么使用 MPBOOTPHYS。

答: 我认为在 boot.S 中还没有页面映射, 因此

Exercise 3.

Task. 修改 kern/pmap.c 中函数 mem_init_mp() (栈的映射)

```
1  总共有 NCPU(8) 个处理器, 对于每个处理器, 分配 KSTKSIZE (8*pgsize) 作为栈, 再做
   KSTKGAP (8*pgsize) 作为保护页,
2  内核栈权限是内核读写以及用户无, 因此权限位设置为 PTE_P|PTE_W
3  static void
4  mem_init_mp(void)
5  {
6      uintptr_t address = KSTACKTOP - KSTKSIZE;
7      for (int i = 0; i < NCPU; i++){
8
9          boot_map_region(kern_pgdir, address, KSTKSIZE, PADDR(percpu_kstacks[i]), PTE_P|PTE_W);
10         address -= (KSTKSIZE + KSTKGAP);
11     }
```

Exercise 4.

Task. 修改 kern/trap.c 中的函数 trap_init_percpu 使得它对所有处理器都有效。(提示 不能使用 ts)

```
1  初始化 每个 CPU 的任务状态段。
2  在原先的基础上修改, 保证对于每个 CPU 都是有效的。(细节比较多)
3  1. thiscpu 指向当前 CPU 对应结构体, 从结构体可以获取 CPU 的 id (id 也可以用 cpunum() 函数获得) 和 ts。
4  2. 将原先的 ts 替换为 thiscpu->cpu_ts 保证每个 CPU 有独立的 ts。
5  3. 修改 esp0 时, 原先对应栈即 KSTACKTOP, 现在不同 CPU 内核栈地址不同, 可通过 Exercise 3 简单计算。
6  4. 修改 gdt 表, 不同 CPU 应放入 不同表项中, 用 (GD_TSS0 >> 3) + thiscpu->cpu_id 计算。
7  5. 加载 TSS 选择器。
8  // Initialize and load the per-CPU TSS and IDT
9  void
10 trap_init_percpu(void)
11 {
12     cprintf("%x\n", thiscpu->cpu_ts);
```

```

13     thiscpu->cpu_ts.ts_esp0 = KSTACKTOP - thiscpu->cpu_id * (KSTKSIZE +
KSTKGAP);
14     thiscpu->cpu_ts.ts_ss0 = GD_KD;
15     thiscpu->cpu_ts.ts_iomb = sizeof(struct Taskstate);
16     gdt[(GD_TSS0 >> 3) + thiscpu->cpu_id] = SEG16(STS_T32A, (uint32_t)
(&thiscpu->cpu_ts),
17         sizeof(struct Taskstate) - 1, 0);
18     gdt[(GD_TSS0 >> 3) + thiscpu->cpu_id].sd_s = 0;
19     ltr(GD_TSS0 + thiscpu->cpu_id * 8);
20     lidt(&idt_pd);
21 }

```

Exercise 5.

Task. 通过 *big kernel lock* 保证所有环境至多只有一个在内核模式中运行。适当的加锁和解锁保证实现。

- 1 1. 在函数 `i386_init()`，在 BSP 唤醒其他 CPU 前加锁，该操作在 `boot_aps` 中执行，因此在该语句前加锁。
- 2 2. 在函数 `mp_main()`，在建立环境后，运行环境前加锁，同时加入运行环境语句。
- 3 3. 在函数 `trap()`，由用户态陷入内核态加锁。
- 4 4. 在函数 `env_run()`，这里应该是切换环境后，和在恢复现场前（这应该在用户模式下执行）。

Question. 为什么同一个时刻只有一个环境在内核模式下，还要对于不同 CPU 分配不同的内核栈而不能共享同一个栈。

答：因为虽然用 *big kernel lock* 保证所有环境中只有一个在内核模式下运行，但是当中断发生的时候，在用户模式下，还未判断是否有锁，硬件部分会将部分寄存器 push 到内核栈中。如果共享栈，可能虽然没有进入内核，但仍然修改了内核栈，导致当前内核程序错误。

Exercise 6.

Task. 轮询调度。

1. 实现 `kern/sched.c` 中函数 `sched_yield()` 且满足一下要求 从上一次推出的环境开始找一个状态为 `ENV_RUNNABLE` 的环境并调用。
2. 修改 `syscall` 实现 `sys_yield()`。
3. 在 `mp_main` 最后添加 `sched_yield()`

4. 建立多个环境运行 user/yield.c (注 user_yield)

```
1  thiscpu->cpu_env 指向当前运行环境（可能不存在，需要特殊处理）
2
3  void
4  sched_yield(void)
5  {
6      struct Env *idle;
7      struct Env*current_env = thiscpu->cpu_env;
8      size_t id = 0;
9      if (current_env != NULL)id = (ENVX(current_env->env_id) + 1) % NENV;
10     for (int i = 0; i < NENV; i++,id = (id + 1) % NENV){
11         if (envs[id].env_status == ENV_RUNNABLE){
12             envs[id].env_cpunum = cpunum();
13             env_run(&envs[id]);
14             return;
15         }
16     }
17     if (current_env != NULL && current_env->env_status == ENV_RUNNING){
18         current_env->env_cpunum = cpunum();
19         env_run(current_env);
20         return;
21     }
22     // sched_halt never returns
23     sched_halt();
24 }
```

Question.

3. 在函数 env_run() 中调用 lr3 切换页目录前后为什么 e 指针保持不变。

答：因为 e 是一个指向环境的指针，该结构维护在内核中，因此所有的环境对这段空间的映射相同。

4. 为什么在切换环境的时候要保存旧环境的寄存器等值，且在什么地方保存。

答：因为如果要恢复环境就需要将寄存器的值恢复，如果不保存则就无法恢复。是在 trapentry.S 的 _alltraps 中，详细可见 lab3。

Exercise 7.

Task. 补充 syscall.c 使完成 fork 操作，通过 user/dumbfork。

```
1 sys_exofork:创建一个新环境，如果成功返回新环境的 ID，失败返回 error_code。
2 这里处理两种 error_code
3 E_NO_FREE_ENV 即没有可用的新的环境，在 env_alloc 函数中若无法分配则会返回该
  error_code。
4 E_NO_MEM 通用在 env_alloc 函数中若无物理页分配时返回该 error_code。
5
6 对于新建的环境，需要以下部分，
7 1. env_status 修改为 ENV_NOT_RUNNABLE
8 2. 寄存器 继承 当前运行环境。
9 3. 新的环境中返回值为 0，即将寄存器 eax 修改为 0。
10 static env_id_t
11 sys_exofork(void)
12 {
13     struct Env *e;
14     int err = env_alloc(&e, curenv->env_id);
15     if (err < 0) return err;
16     e->env_status = ENV_NOT_RUNNABLE;
17     e->env_tf = curenv->env_tf;
18     e->env_tf.tf_regs.reg_eax = 0;
19     return e->env_id;
20     panic("sys_exofork not implemented");
21 }
```

```
1 sys_env_set_status:修改当前环境状态，且只能修改为 ENV_RUNNABLE 或
  ENV_NOT_RUNNABLE。
2 函数 env_id2env 将 env_id 转换为对应 env 结构体。
3 需要处理 error_code
4 E_BAD_ENV
5 1. 在 env_id2env 中若非法的 env 会返回该 error_code
6 2. 权限不足，在调用 env_id2env 中将第三个参数设置为 1 即会检查权限。
7 E_INVALID 要求 status 为 ENV_RUNNABLE 或 ENV_NOT_RUNNABLE
8
9 修改 env 对应状态即可。
10 static int
11 sys_env_set_status(env_id_t env_id, int status)
12 {
13     if (status != ENV_RUNNABLE && status != ENV_NOT_RUNNABLE)
14         return -E_INVALID;
```

```

15     struct Env *e;
16     int err = envid2env(envid,&e,1);
17     if (err < 0)return err;
18     e->env_status = status;
19     return 0;
20     panic("sys_env_set_status not implemented");
21 }

```

```

1  为环境分配分配页面
2  参数限定：
3  perm 中 PTE_U 和 PTE_P 位为 1 ， 且 PTE_AVAIL 和 PTE_W 可以为 0 或 1，其他位
   必须为 0 。
4  错误处理：
5  E_BAD_ENV:判读环境是否合法
6  E_INVAL:判读虚拟地址 1. 对齐, 2.不能超过UTOP 3.权限位合理
7  E_NO_MEM:能够建立映射（先用 page_alloc 从系统分配物理页，再用 page_insert 建立
   虚拟地址和物理页之间映射）
8
9  static int
10 sys_page_alloc(envid_t envid, void *va, int perm)
11 {
12     // Hint: This function is a wrapper around page_alloc() and
13     //     page_insert() from kern/pmap.c.
14     //     Most of the new code you write should be to check the
15     //     parameters for correctness.
16     //     If page_insert() fails, remember to free the page you
17     //     allocated!
18
19     // LAB 4: Your code here.
20     struct Env *e;
21     int err = envid2env(envid,&e,1);
22     if (err < 0)return err;
23     if ((uint32_t)va >= UTOP || PGOFF(va))return -E_INVAL;
24     if ((perm & ~(PTE_AVAIL|PTE_W))^(PTE_U|PTE_P))
25         return -E_INVAL;
26     struct PageInfo *page = page_alloc(ALLOC_ZERO);
27     if (page == NULL) return -E_NO_MEM;
28     err = page_insert(e->env_pgdir,page,va,perm);
29     if (err<0){
30         page_free(page);
31         return -E_NO_MEM;
32     }

```

```

33     return 0;
34     panic("sys_page_alloc not implemented");
35 }

```

```

1  建立环境之间的页面映射关系。
2  参数限定：
3  perm：首先和 page_alloc 要求相同，在基础上还要求对于 只读页 不能加上 可写的参
    数。
4
5  错误处理：
6  E_BAD_ENV：环境错误，这里有源环境和目的环境，要求两个都合法。
7  E_INVALID：
8  1. 源地址和目的地址 对齐且不超过 UTOP
9  2. 源环境中对应的地址映射。
10 3. 权限位正确，包括 （要求对于 只读页 不能加上 可写的参数）。
11 4. 能够在目的环境中分配物理地址并建立映射。
12 static int
13 sys_page_map(envid_t srcenvid, void *srcva,
14               envid_t dstenvid, void *dstva, int perm)
15 {
16     struct Env*esrc,*edst;
17     int errsrc = envid2env(srcenvid,&esrc,1),errdst =
    envid2env(dstenvid,&edst,1);
18     if (errsrc < 0 || errdst < 0)return -E_BAD_ENV;
19     if ((uint32_t)srcva >= UTOP || PGOFF(srcva) || (uint32_t)dstva >=
    UTOP || PGOFF(dstva))return -E_INVALID;
20     pte_t* pte;
21     struct PageInfo* page = page_lookup(esrc->env_pgdir, srcva, &pte);
22     if (page == NULL) return -E_INVALID;
23     if ((perm & ~(PTE_AVAIL|PTE_W))^(PTE_U|PTE_P))
24         return -E_INVALID;
25     if ((perm & PTE_W)&&!(pte & PTE_W))return -E_INVALID;
26     struct PageInfo* pagedst = page_alloc(ALLOC_ZERO);
27     if (page == NULL) return -E_NO_MEM;
28     int err = page_insert(edst->env_pgdir,page,dstva,perm);
29     if (err < 0){
30         page_free(pagedst);
31         return -E_NO_MEM;
32     }
33     return 0;
34     panic("sys_page_map not implemented");
35 }

```



```

1  取消映射
2  错误处理：
3  E_BAD_ENV 环境错误。
4  E_INVAL 地址对齐和 不超过 UTOP
5
6  用 page_remove 删除。
7  static int
8  sys_page_unmap(envid_t envid, void *va)
9  {
10     struct Env *e;
11     int err = envid2env(envid,&e,1);
12     if (err < 0)return err;
13     pte_t*pte;
14     struct PageInfo* page = page_lookup(e->env_pgdir,va,&pte);
15     if (pte == NULL || !(*pte & PTE_W))return -E_BAD_ENV;
16     if ((uint32_t)va >= UTOP || PGOFF(va)) return -E_INVAL;
17     page_remove(e->env_pgdir,va);
18     return 0;
19     panic("sys_page_unmap not implemented");
20 }

```

Part B: Copy-on-Write Fork

Exercise 8.

Task. 不同段的 page fault 处理过程并不相同，因此需要用 *page fault handler entrypoint* 记录。实现 `sys_env_set_pgfault_upcall` 函数。

```

1  修改 对应的 env_pgfault_upcall 项即可。
2  static int
3  sys_env_set_pgfault_upcall(envid_t envid, void *func)
4  {
5     struct Env *e;
6     int err = envid2env(envid,&e,1);
7     if (err < 0)return err;
8     e->env_pgfault_upcall = func;
9     return 0;
10 }

```

Exercise 9.

Task. 实现 kern/trap.c 中 page_fault_handler 对于用户程序的 page fault 的处理。

```
1  这里 处理 在用户环境中 发生的page_fault
2  主要流程包括以下几个步骤
3  1. 找对应处理函数。
4  2. 在异常栈中分配空间保存现场。
5  3. 修改环境 eip 和 esp 并运行处理函数。
6
7  第一步中, 在环境中查找 page_fault_upcall , 如果有会在
   sys_env_set_pgfault_upcall 建立映射。
8  第二步中, 在用户异常栈中分配栈帧用于保存现场,
9  1. 首先要分配从外部跳入到 page_fault 还是 从 page_fault 处理过程中 再次跳入
   page_fault。这里检查原现场的 esp 即可
10 (应为递归调用是 esp 一定在用户异常栈之间) 。
11 如果从外部跳入, 则从异常栈头开始, 否则 从上一次 -1 word 开始, (具体为什么要多分配
   4 字节空间, 是因为在回退的时候需要保存地址, 在第 Exercise 10 中就可发现)
12 2.用 user_mem_assert 去检查分配的栈帧 是否是 用户可读可写的,user_mem_check 主
   要检查越界和权限两方面。
13 3.保存现场, 这里仿照lab 中的图赋值即可。
14 第三步中, 将指令跳转到处理程序中, 栈切换位用户异常栈, 并运行。
15 Hint 中 tf 即当前环境的 env_tf
16
17 void
18 page_fault_handler(struct Trapframe *tf)
19 {
20     uint32_t fault_va;
21     fault_va = rcr2();
22     if ( (tf->tf_cs&1)!=1 )
23         panic("page_fault_handler: kernel page fault!\n");
24
25     // LAB 4: Your code here.
26     if (curenv->env_pgfault_upcall!=NULL){
27         // cprintf("%x\n",tf->tf_esp);
28         struct UTrapframe *utf = (tf->tf_esp >= UXSTACKTOP || tf->tf_esp <
   UXSTACKTOP - PGSIZE) ?
29         (struct UTrapframe *) (UXSTACKTOP - sizeof(struct UTrapframe)) :
   (struct UTrapframe *) (tf->tf_esp - 4 - sizeof(struct UTrapframe));
30         // cprintf("find %x\n",utf);
31         user_mem_assert(curenv,(const void*)utf,sizeof(struct
   UTrapframe),PTE_U|PTE_W|PTE_P);
```

```

32     // cprintf("find2\n");
33     utf->utf_esp = tf->tf_esp;
34     utf->utf_eflags = tf->tf_eflags;
35     utf->utf_eip = tf->tf_eip;
36     utf->utf_regs = tf->tf_regs;
37     utf->utf_err = tf->tf_trapno;
38     utf->utf_fault_va = fault_va;
39
40     tf->tf_eip = (uintptr_t)curenv->env_pgfault_upcall;
41     // cprintf("eip %x\n",tf->tf_eip);
42     tf->tf_esp = (uintptr_t) utf;
43     env_run(curenv);
44 }
45 // Destroy the environment that caused the fault.
46 cprintf("[%08x] user fault va %08x ip %08x\n",
47     curenv->env_id, fault_va, tf->tf_eip);
48 print_trapframe(tf);
49 env_destroy(curenv);
50 }
51

```

Exercise 10.

Task. 用汇编实现 page_fault handler 回掉程序。

```

1  这里主要实现从处理函数结束恢复现场的过程。
2  观察栈帧结构可以知道，UTrapframe 如此设计的妙处，
3  即先恢复寄存器，再恢复标志寄存器，最后恢复 esp。
4
5  主要想法，处理函数的调用实际上是从一个栈跳转到了异常处理栈中并处理过程，在恢复的时
   候，我们直接从异常处理栈中 ret 不能保证寄存器完整，因此我们先在原栈中构建一个新的小
   的栈帧，再在原栈中 ret。
6  1.我们先把 eip 写到原栈的下面，仿佛是在原栈中进行因此内部的过程调用，并保存的 eip。
7  2.对齐恢复所有的常用寄存器。
8  3.这里先跳过 eip，再恢复标志寄存器。
9  4.最后由于栈帧的最顶部保存的是原现场的 esp，因此用 pop %esp 指令可以得假装过程调用
   的栈地址。
10 5.再用 ret 恢复到真实 eip 和 esp
11 .text
12 .globl _pgfault_upcall

```

```

13 _pgfault_upcall:
14     // Call the C page fault handler.
15     pushl %esp          // function argument: pointer to UTF
16     movl _pgfault_handler, %eax
17     call *%eax
18     addl $4, %esp       // pop function argument
19
20     movl 0x28(%esp), %eax
21     subl $0x4, 0x30(%esp)
22     movl 0x30(%esp), %ebx
23     movl %eax, (%ebx)
24     addl $8, %esp
25     popal
26
27     addl $4, %esp
28     popfl
29
30     mov (%esp), %esp
31
32     ret

```

Exercise 11.

Task. 完善 set_pgfault_handler 函数

```

1  第一次处理的时候要 建立异常栈 和 将指针连接到 环境处理函数中，
2  之后修改处理函数只需要修改函数指针即可。
3
4  void
5  set_pgfault_handler(void (*handler)(struct UTrapframe *utf))
6  {
7      int r;
8
9      if (_pgfault_handler == 0) {
10         envid_t eid = sys_getenvid();
11         r = sys_page_alloc(eid, (void *) (UXSTACKTOP - PGSIZE),
PTE_P|PTE_U|PTE_W);
12         if (r<0) panic("set_pgfault_handler: sys_page_alloc\n");
13         r = sys_env_set_pgfault_upcall(eid, _pgfault_upcall);
14         if (r<0) panic("set_pgfault_handler: sys_env_set_pgfault_upcall\n");

```

```

15     }
16
17     // Save handler pointer for assembly to call.
18     _pgfault_handler = handler;
19 }

```

测试到过程中还发现到一个 lab3 中的bug

```

1  lab3 的 bug
2  应先检查 pte 是否为空, 否则对pte 引用会导致内核的 page_fault。
3  user_mem_check(struct Env *env, const void *va, size_t len, int perm)
4  {
5      ...
6      for (uintptr_t address = start; address < end; address+= PGSIZE){
7          if (address >= ULIM){
8              user_mem_check_addr = (address>(uintptr_t)va)?address:
9              (uintptr_t)va;
10             return -E_FAULT;
11         }
12         pte_t *pte = pgdir_walk(env->env_pgdir, (void *)address, 0);
13         if (pte == NULL || (*pte & perm) != perm){
14             user_mem_check_addr = (address>(uintptr_t)va)?address:
15             (uintptr_t)va;
16             return -E_FAULT;
17         }
18     }
19     return 0;
20 }

```

Exercise 12.

Task. 完善 fork 函数。

```

1  从当前环境将 pn 对应页映射到 env 环境中,
2  1. pn * pgsize 即 该页的虚拟地址,
3  2. perm 首先设置为 PTE_U 和 PTE_P, 如果 是 可写或写时 复制则同时加上 PTE_COW
4  3. 映射到 child env 中
5  4. 若当前环境该页 非 写时复制 则再回映射到 parent env 中。
6  static int
7  duppage(env_t env, unsigned pn)

```

```

8 {
9     int r;
10    // LAB 4: Your code here. envid_t myenvid = sys_getenvid();
11    void *va = (void*)(pn * PGSIZE);
12    pte_t pte = uvpt[pn];
13    envid_t envid_parent = sys_getenvid();
14    int perm = PTE_U|PTE_P|(((pte&(PTE_W|PTE_COW))>0)?PTE_COW:0);
15    int err;
16    err = sys_page_map(envid_parent, va, envid, va, perm);
17    if (err < 0) return err;
18    if ((perm|~pte)&PTE_COW){
19        err = sys_page_map(envid_parent, va, envid_parent, va, perm);
20        if (err < 0) return err;
21    }
22    return 0;
23    panic("duppage not implemented");
24    return 0;
25 }

```

1 错误处理函数

2 1. 首先检查 fault is a write 和 当前页时 PTE_COW

3 2. 在 PFTEMP 分配一个用户可读可写的页 , 然后复制, 并建立映射。

```

4 static void
5 pgfault(struct UTrapframe *utf)
6 {
7     void *addr = (void *) utf->utf_fault_va;
8     uint32_t err = utf->utf_err;
9     int r;
10    if ((err&FEC_WR)==0)
11        panic("pgfault: error!\n");
12    if ((uvpt[PGNUM(addr)]&PTE_COW)==0)
13        panic("pgfault: error!\n");
14    envid_t envid = sys_getenvid();
15    int err = sys_page_alloc(envid, (void*)PFTEMP, PTE_P|PTE_U|PTE_W);
16    if (err<0) panic("pgfault: error!\n");
17    addr = ROUNDDOWN(addr, PGSIZE);
18    memcpy(PFTEMP, addr, PGSIZE);
19    err = sys_page_map(envid,
20    (void*)PFTEMP, envid, addr, PTE_P|PTE_U|PTE_W);
21    if (err<0) panic("pgfault: error!\n");
22    err = sys_page_unmap(envid, PFTEMP);
23    if (err<0) panic("pgfault: error!\n");

```

```

23     return;
24     panic("pgfault not implemented");
25 }

```

```

1  fork 函数
2  函数流程
3  1. parent: 将 page_fault_handle 设置为之前定义的 pgfault().
4  2. parent: 调用 sys_exofork 建立新的 环境。
5  3.1 parent: 对于 UTOP 以下的页, 调用之前定义的 duppage 将标记为 w 或者 cow 的
    页映射到子进程的地址空间, 标记为 cow , 同时父进程本身也修改为 cow。(先子后父)。异
    常栈则是需要申请一个新的页。
6  3.2 child: 准备运行状态
7  4 parent 修改 child 状态为 runnable
8
9  envid_t
10 fork(void)
11 {
12     // LAB 4: Your code here.
13     set_pgfault_handler(pgfault);
14     envid_t envid = sys_exofork();
15     if (envid<0)
16         panic("fork : error!\n");
17     if (envid==0){
18         thisenv = envs+ENVX(sys_getenvid());
19         return envid;
20     }
21     envid_t envid_child = envid;
22     envid_t envid_parent = sys_getenvid();;
23     for (uintptr_t addr=0;addr<USTACKTOP;addr+=PGSIZE)
24         if ((uvpd[PDX(addr)]&PTE_P)&&(uvpt[PGNUM(addr)]&PTE_P))
25             duppage(envid_child,PGNUM(addr));
26     if (sys_page_alloc(envid_child,(void *) (UXSTACKTOP-
    PGSIZE),PTE_U|PTE_W|PTE_P)<0)panic("fork : error!\n");
27     if (sys_env_set_pgfault_upcall(envid_child,(envs +
    ENVX(envid_parent))->env_pgfault_upcall)<0)panic("fork : error!\n");
28     if (sys_env_set_status(envid_child, ENV_RUNNABLE)<0)panic("fork :
    error!\n");
29     return envid_child;
30     panic("fork not implemented");
31 }

```

Part C: Preemptive Multitasking and Inter-Process communication (IPC)

Exercise 13.

Task. 修改 kern/trapentry.S 和 kern/trap.c 为 IDT 新增关于 IRQ 的表项, 修改 kern/env.c 中函数 env_alloc() 保证用户环境能够中断。

1 | 类似 lab3 为异常增加表项。

Exercise 14.

Task. 实现时钟中断时实现切换进程。

```
1 | 这里只处理 时钟中断, lapic_eoi 确认中断, 并切换调度。
2 | static void
3 | trap_dispatch(struct Trapframe *tf)
4 | {
5 |     ...
6 |     if (tf->tf_trapno == IRQ_OFFSET + IRQ_SPURIOUS) {
7 |         cprintf("Spurious interrupt on irq 7\n");
8 |         print_trapframe(tf);
9 |         return;
10 |    }
11 |    switch (tf->tf_trapno){
12 |        case IRQ_OFFSET + IRQ_TIMER: {
13 |            lapic_eoi();
14 |            sched_yield();
15 |            return;
16 |        }
17 |    }
18 |    ...
19 | }
```


Exercise 15.

Task.

```
1  进程间信息通信发送信息,
2  错误处理
3  E_BAD_ENV: 环境错误
4  E_IPC_NOT_RECV: 目标进程不允许通信, 用 env_ipc_recving 标记
5  E_INVALID: 对于 srcva < UTOP
6  首先要 页对齐 且权限位正确
7  发送环境中要有该页面映射
8  不能对只读页面进行可写权限的映射。
9  E_NO_MEM: 缺乏物理页
10
11 否则 修改 ipc_recving, ipc_from, ipc_value, ipc_perm, status 为对应的值。
12
13 static int
14 sys_ipc_try_send(envid_t envid, uint32_t value, void *srcva, unsigned
    perm)
15 {
16     // LAB 4: Your code here.
17     envid_t src_envid = sys_getenvid();
18     struct Env *e;
19     int err;
20     err = envid2env(envid, &e, 0);
21     // cprintf("err %x\n", err);
22     if (err < 0) return err;
23     if (!e->env_ipc_recving) return -E_IPC_NOT_RECV;
24     if ((uint32_t)srcva < UTOP && PGOFF(srcva)) return -E_INVALID;
25     if ((perm & ~(PTE_AVAIL|PTE_W)) ^ (PTE_U|PTE_P)) return -E_INVALID;
26     pte_t *pte;
27     struct PageInfo *page = page_lookup(curenv->env_pgdir, srcva, &pte);
28     if ((uint32_t)srcva < UTOP && page == NULL) return -E_INVALID;
29     if ((uint32_t)srcva < UTOP && (perm & PTE_W) && (~*pte &
        PTE_W)) return -E_INVALID;
30     if ((uint32_t)srcva < UTOP){
31         err = page_insert(e->env_pgdir, page, e->env_ipc_dstva, perm);
32         if (err < 0) return err;
33     }
34     // cprintf("find\n");
35     e->env_ipc_recving = false;
36     e->env_ipc_from = src_envid;
```

```

37     e->env_ipc_value = value;
38     e->env_ipc_perm = ((uint32_t)srcva < UTOP)?perm:0;
39     e->env_status = ENV_RUNNABLE;
40
41     e->env_tf.tf_regs.reg_eax = 0;
42     return 0;
43     panic("sys_ipc_try_send not implemented");
44 }

```

1 进程接受信息，如果 dstva < UTOP 说明要页映射，因此要页对齐。
2 将 ipc_recving 修改为true，修改dstva，并将进程设置成 不能再被调用。

```

3 static int
4 sys_ipc_recv(void *dstva)
5 {
6     // LAB 4: Your code here.
7     if ((uintptr_t)dstva<UTOP&&PGOFF(dstva)!=0)return -E_INVAL;
8     envid_t envid = sys_getenvid();
9     struct Env *e;
10    e->env_ipc_recving = true;
11    e->env_ipc_dstva = dstva;
12    e->env_status = ENV_NOT_RUNNABLE;
13    sys_yield();
14    return 0;
15 }

```

1 若 pg 为 NULL 则修改为 UTOP 表示不进行页面传递。
2 若接收成功则保存原环境和权限，并返回value 否则清空，返回 error_code。

```

3
4 int32_t
5 ipc_recv(envid_t *from_env_store, void *pg, int *perm_store)
6 {
7     // LAB 4: Your code here.
8
9     if (pg==NULL)pg=(void *)UTOP;
10    // cprintf("now %x\n",pg);
11    int err = sys_ipc_recv(pg);
12    // cprintf("%x\n",err);
13    if (err < 0){
14        if (from_env_store != NULL)*from_env_store=0;
15        if (perm_store != NULL)*perm_store=0;

```

```

16     return err;
17 }else {
18     if (from_env_store != NULL) *from_env_store = thisenv->env_ipc_from;
19     if (perm_store != NULL) *perm_store = thisenv->env_ipc_perm;
20     return thisenv->env_ipc_value;
21 }
22 panic("ipc_recv not implemented");
23 return 0;
24 }

```

```

1 一直发送，直至成功或出错。
2 void
3 ipc_send(envid_t to_env, uint32_t val, void *pg, int perm)
4 {
5     // LAB 4: Your code here.
6     if (pg==NULL)pg=(void *)UTOP;
7     int err;
8     while (true){
9         err = sys_ipc_try_send(to_env,val,pg,perm);
10        // cprintf("%x\n",err);
11        if (err == -E_IPC_NOT_RECV)sys_yield();
12        else if (err < 0)panic("ipc send: error!\n");
13        else return;
14    }
15    return;
16    panic("ipc_send not implemented");
17 }

```

Challenge

Challenge! Implement a shared-memory `fork()` called `sfork()`. This version should have the parent and child *share* all their memory pages (so writes in one environment appear in the other) except for pages in the stack area, which should be treated in the usual copy-on-write manner. Modify `user/forktree.c` to use `sfork()` instead of regular `fork()`. Also, once you have finished implementing IPC in part C, use your `sfork()` to run `user/pingpongs`. You will have to find a new way to provide the functionality of the global `thisenv` pointer.

Task. 实现 `sfork` 使得对于所有用户空间 除栈外 共享。

```
1  难点, thisenv 因为 thisenv 位于用户空间中, 父子进程共享用户空间, 因此父子进程的
   thisenv 会相同 (即如果再子进程修改 thisenv, 那么父进程也会被修改导致错误。)
2  考虑到用户栈是会被共享的, 因此将 thisenv 保存在用户栈中。
3
4  原先 thisenv 为全局变量指向当前 env 的指针, 现修改为 pthisenv 为指向 当前 env
   的指针的指针。
5  在函数局部定义 local_thisenv 为指向当前 env 的指针, 由于在局部定义, 因此位于用户
   栈中。
6  将 pthisenv 指向 local_thisenv, 并将 thisenv define 为 (*pthisenv) (该操作
   只是为了保证最小化修改, 只需要在库文件中加入该宏定义, 否则对于原代码中使用到 thisenv
   的地方修改也可)
7
8  #include <inc/lib.h>
9  extern void umain(int argc, char **argv);
10 const char *binaryname = "<unknown>";
11 #define thisenv (*pthisenv)
12 const volatile struct Env **pthisenv;
13
14 void
15 libmain(int argc, char **argv)
16 {
17     const volatile struct Env *local_thisenv = (envs +
18     ENVX(sys_getenvvid()));
19     pthisenv = &local_thisenv;
20     if (argc > 0)
21         binaryname = argv[0];
22     umain(argc, argv);
23     exit();
24 }
```

```
1  类似于 fork, 在映射过程中有所修改, 对于除栈外空间用 sys_page_map 从父进程映射到子
   进程, 对于 用户栈 还是用 cow 方式。
2  int
3  sfork(void)
4  {
5      set_pgfault_handler(pgfault);
6      envvid_t envvid = sys_exofork();
7      cprintf("envvid :%x\n", envvid);
8      if (envvid < 0)
9          panic("sfork : error!\n");
10 }
```

```

10  if (envid==0){
11      thisenv = envs+ENVX(sys_getenvid());
12      return envid;
13  }
14  envid_t envid_child = envid;
15  envid_t envid_parent = sys_getenvid();
16  for (uintptr_t addr=UTEXT;addr<USTACKTOP - PGSIZE;addr+=PGSIZE)
17      if ((uvpd[PDX(addr)]&PTE_P)&&(uvpt[PGNUM(addr)]&PTE_P)){
18          // if ((uvpd[PDX(addr)]&PTE_P)&&(uvpt[PGNUM(addr)]&PTE_P))
19          //      duppage(envid_child,PGNUM(addr));
20          sys_page_map(envid_parent,(void*)addr,envid_child,
(void*)addr,PTE_U|PTE_P|PTE_W);
21      }
22      duppage(envid_child,PGNUM(USTACKTOP - PGSIZE));
23      if (sys_page_alloc(envid_child,(void *) (UXSTACKTOP-
PGSIZE),PTE_U|PTE_W|PTE_P)<0)panic("sfork : error!\n");
24      if (sys_env_set_pgfault_upcall(envid_child,(envs +
ENVX(envid_parent))->env_pgfault_upcall)<0)panic("sfork : error!\n");
25      if (sys_env_set_status(envid_child, ENV_RUNNABLE)<0)panic("sfork :
error!\n");
26      return envid_child;
27
28  panic("sfork not implemented");
29  return -E_INVALID;
30  }

```

测试:

```
0
SMP: CPU 0 found 1 CPU(s)
enabled interrupts: 1 2
[00000000] new env 00001000
eebdfde4
[00001000] new env 00001001
envid :1001
i am 00001000; thisenv is 0xeec00000
envid :0
send 0 from 1000 to 1001
1001 got 0 from 1000 (thisenv is 0xeec0007c 1001)
1000 got 1 from 1001 (thisenv is 0xeec00000 1000)
1001 got 2 from 1000 (thisenv is 0xeec0007c 1001)
1000 got 3 from 1001 (thisenv is 0xeec00000 1000)
1001 got 4 from 1000 (thisenv is 0xeec0007c 1001)
1000 got 5 from 1001 (thisenv is 0xeec00000 1000)
1001 got 6 from 1000 (thisenv is 0xeec0007c 1001)
1000 got 7 from 1001 (thisenv is 0xeec00000 1000)
1001 got 8 from 1000 (thisenv is 0xeec0007c 1001)
1000 got 9 from 1001 (thisenv is 0xeec00000 1000)
[00001000] exiting gracefully
[00001000] free env 00001000
1001 got 10 from 1000 (thisenv is 0xeec0007c 1001)
[00001001] exiting gracefully
[00001001] free env 00001001
No runnable environments in the system!
```

运行 pingpongs.c 的结果，可以发现 val 在父进程和子进程之间得到共享。