

---

# Lab #4

## 陈威宇

### Exercise 1

mmio\_map\_region 将 pa 开始的 size 个字节 map 到 MMIO region.

```
1 void *
2 mmio_map_region(physaddr_t pa, size_t size)
3 {
4     static uintptr_t base = MMIIOBASE;
5
6     int mo = pa%PGSIZE;
7     size+=mo;
8     pa-=mo;
9     if (size%PGSIZE!=0)
10         size=size-size%PGSIZE+PGSIZE;
11     if (base + size>MMIOBASE+0x400000)
12         panic("MMIO region is not enough!!!");
13     boot_map_region(kern_pgdir, base, size, pa, PTE_PCD|PTE_PWT|PTE_W);
14     uint32_t t = base + mo;
15     base += size;
16     return (void *)t;
17 }
```

### Exercise 2

修改了 page\_init 使得 MPENTRY\_PADDR 开始的一个页不是 free page.

```
1 void
2 page_init(void)
3 {
4     page_free_list=NULL;
5     size_t i;
6     for (i = 0; i < npages; i++) {
7         pages[i].pp_ref = 0;
8         if (i==0 || (i*PGSIZE >= IOPHYSMEM && i*PGSIZE < EXTPHYSMEM) || (i*PGSIZE>=0x100000 && i*P
9             pages[i].pp_link = NULL;
10     }
11     else{
12         pages[i].pp_link = page_free_list;
```

---

```
13     page_free_list = &pages[i];
14 }
15 }
16 }
17
```

## Questions after Exercise 2

1

因为 mpentry.S 是在高地址链接的 (KERNBASE 以上), 但是被 bootaps load 到低地址 (MPENTRY\_ADDR), 所以要用 MPBOOTPHYS 求物理地址.

## Exercise 3

mem\_init\_mp 给每个 CPU 都 map 了一个 kernel stack.

```
1  static void
2  mem_init_mp(void)
3  {
4      for (int i=0; i<NCPU; ++i){
5          uintptr_t kstacktop_i = KSTACKTOP - i * (KSTKSIZE + KSTKGAP);
6          boot_map_region(kern_pgdir, kstacktop_i-KSTKSIZE, KSTKSIZE, PADDR(percpu_kstacks[i]), PTE_
7      }
8  }
```

## Exercise 4

trap\_init\_percpu 对每个 CPU 初始化了 TSS and IDT.

```
1  void
2  trap_init_percpu(void)
3  {
4      /*
5       ts.ts_esp0 = KSTACKTOP;
6       ts.ts_ss0 = GD_KD;
7       ts.ts_iomb = sizeof(struct Taskstate);
8      */
9      (thiscpu->cpu_ts.ts_esp0) = KSTACKTOP - (thiscpu->cpu_id) * (KSTKSIZE + KSTKGAP);
10     (thiscpu->cpu_ts.ts_ss0) = GD_KD;
11     (thiscpu->cpu_ts.ts_iomb) = sizeof(struct Taskstate);
12     // Initialize the TSS slot of the gdt.
```

---

```

13  /*
14     gdt[GD_TSS0 >> 3] = SEG16(STS_T32A, (uint32_t) (&ts),
15         sizeof(struct Taskstate) - 1, 0);
16     gdt[GD_TSS0 >> 3].sd_s = 0;
17  */
18     gdt[(GD_TSS0 >> 3)+cpunum()] = SEG16(STS_T32A, (uint32_t) (&(thiscpu->cpu_ts)),
19         sizeof(struct Taskstate) - 1, 0);
20     gdt[(GD_TSS0 >> 3)+cpunum()].sd_s = 0;
21
22     // Load the TSS selector (like other segment selectors, the
23     // bottom three bits are special; we leave them 0)
24     ltr(GD_TSS0 + ((cpunum())<<3));
25
26     // Load the IDT
27     lidt(&idt_pd);
28 }

```

## Exercise 5

按文档指示在对应位置加上 kernel lock 的语句就行了.

## Questions after Exercise 5

2

因为发生异常或中断时, 硬件会把寄存器状态放到 kernel stack 上, 这是不需要 kernel lock 的.

## Exercise 6

sched\_yield 实现了 round-robin scheduling.

```

1  void
2  sched_yield(void)
3  {
4      struct Env *idle;
5
6      int id;
7      if (curenv==NULL)
8          id=-1;
9      else
10         id = ENVX(curenv->env_id);

```

---

```

11  int h = id;
12  while (1){
13      ++h;
14      if (h==id)
15          break;
16      if (h>=NENV){
17          if (id<=0)
18              break;
19          h=0;
20      }
21      if (envs[h].env_status == ENV_RUNNABLE){
22          env_run(&envs[h]);
23      }
24
25  }
26  if (id>=0 && envs[id].env_status == ENV_RUNNING)
27      env_run(&envs[id]);
28
29  // sched_halt never returns
30  sched_halt();
31  }

```

## Questions after Exercise 6

3

虽然换了一个 env, 但是它们的虚拟内存的 UTOP 以上部分是一样的, 所以不受影响.

4

需要保存旧的 env 的寄存器等, 为了之后接着运行的时候能运行. 在 trap 函数里的时候存到了 env\_tf, 之后 env\_pop\_tf() 的时候又放回了 kernel stack.

## Exercise 7

sys\_exofork 分配了一个新的 environment.

```

1  static envid_t
2  sys_exofork(void)
3  {
4      struct Env * t;
5      int e = env_alloc(&t, curenv->env_id);

```

---

```

6     if (e<0) return e;
7     (t->env_status) = ENV_NOT_RUNNABLE;
8     (t->env_tf) = curenv->env_tf;
9     t->env_tf.tf_regs.reg_eax = 0;
10    return t->env_id;
11 }

```

sys\_env\_set\_status 给一个 environment 设置了 status.

```

1  static int
2  sys_env_set_status(envid_t envid, int status)
3  {
4      if (status != ENV_RUNNABLE && status != ENV_NOT_RUNNABLE)
5          return -E_INVALID;
6      struct Env * t;
7      int e = env_id2env(envid, &t, 1);
8      if (e<0)
9          return e;
10     t->env_status = status;
11     return 0;
12 }

```

sys\_page\_alloc 分配了一页物理页给一个 env 的 va 开始的虚拟页来映射.

```

1  static int
2  sys_page_alloc(envid_t envid, void *va, int perm)
3  {
4      if ((uint32_t)va >= UTOP || (uint32_t)va%PGSIZE!=0)
5          return -E_INVALID;
6      if ((perm & (PTE_U | PTE_P)) != (PTE_U | PTE_P) || (((perm ^ (PTE_U | PTE_P)) | (PTE_AVAIL
7          return -E_INVALID;
8      struct Env * t;
9      int e = env_id2env(envid, &t, 1);
10     if (e<0)
11         return e;
12     struct PageInfo * pp = page_alloc(1);
13     if (pp == NULL)
14         return -E_NO_MEM;
15     e = page_insert(t->env_pgdir, pp, va, perm);
16     if (e<0)
17         return e;
18     return 0;
19 }

```

---

sys\_page\_map 把一个 page mapping 从一个 env 复制到另一个 env.

```
1 static int
2 sys_page_map(envid_t srcenvid, void *srcva,
3             envid_t dstenvid, void *dstva, int perm)
4 {
5     if ((uint32_t)srcva >= UTOP || (uint32_t)srcva%PGSIZE!=0 || (uint32_t)dstva >= UTOP || (uint32_t)dstva%PGSIZE!=0)
6         return -E_INVALID;
7     if ((perm & (PTE_U | PTE_P)) != (PTE_U | PTE_P) || (((perm ^ (PTE_U | PTE_P)) | (PTE_AVAIL | PTE_W)) & PTE_W))
8         return -E_INVALID;
9     struct Env * srcenv;
10    int e;
11    e = envid2env(srcenvid, &srcenv, 1);
12    if (e<0)
13        return e;
14    struct Env * dstenv;
15    e = envid2env(dstenvid, &dstenv, 1);
16    if (e<0)
17        return e;
18    pte_t * pt;
19    struct PageInfo * pp = page_lookup(srcenv->env_pgdir, srcva, &pt);
20    if (pp == NULL)
21        return -E_INVALID;
22    if (((*pt)&PTE_W)==0 && (perm & PTE_W))
23        return -E_INVALID;
24    e = page_insert(dstenv->env_pgdir, pp, dstva, perm);
25    if (e<0)
26        return e;
27    return 0;
28 }
```

sys\_page\_unmap 把一个 env 的虚拟地址 va 开始的页 unmap 了.

```
1 static int
2 sys_page_unmap(envid_t envid, void *va)
3 {
4     if ((uint32_t)va >= UTOP || (uint32_t)va%PGSIZE != 0)
5         return -E_INVALID;
6     struct Env * env;
7     int e = envid2env(envid, &env, 1);
8     if (e<0)
9         return e;
```

---

```
10  page_remove(env->env_pgdir, va);
11  return 0;
12 }
```

## Exercise 8

sys\_env\_set\_pgfault\_upcall 对一个 env 设置了 page fault upcall.

```
1  static int
2  sys_env_set_pgfault_upcall(envid_t envid, void *func)
3  {
4      // LAB 4: Your code here.
5      struct Env * env;
6      int e = envid2env(envid, &env, 1);
7      if (e<0)
8          return e;
9      env->env_pgfault_upcall = func;
10     return 0;
11 }
```

## Exercise 9

page\_fault\_handler 在 user exception stack 上放了个 user trap frame, 然后去跑这个 env, 从 page fault upcall 的位置开始跑, 让 env 解决 page fault.

```
1  void
2  page_fault_handler(struct Trapframe *tf)
3  {
4      uint32_t fault_va;
5
6      // Read processor's CR2 register to find the faulting address
7      fault_va = rcr2();
8
9      // Handle kernel-mode page faults.
10
11     // LAB 3: Your code here.
12     if ((tf->tf_cs & 3) == 0) {
13         panic("kernel-mode page fault!!!");
14     }
15
16     // LAB 4: Your code here.
```

---

```

17  if ((curenv->env_pgfault_upcall) != NULL){
18      struct UTrapframe *utf;
19
20      if (tf->tf_esp >= UXSTACKTOP - PGSIZE && tf->tf_esp < UXSTACKTOP) {
21          *(uint32_t *) (tf->tf_esp - 4) = 0;
22          utf = (struct UTrapframe *) (tf->tf_esp - 4 - sizeof(struct UTrapframe));
23      } else {
24          utf = (struct UTrapframe *) (UXSTACKTOP - sizeof(struct UTrapframe));
25      }
26
27      user_mem_assert(curenv, (void *)utf, sizeof(struct UTrapframe), PTE_W | PTE_U);
28
29      utf->utf_esp = tf->tf_esp;
30      utf->utf_eflags = tf->tf_eflags;
31      utf->utf_eip = tf->tf_eip;
32      utf->utf_regs = tf->tf_regs;
33      utf->utf_err = tf->tf_err;
34      utf->utf_fault_va = fault_va;
35
36      tf->tf_esp = (uint32_t)utf;
37      tf->tf_eip = (uint32_t)curenv->env_pgfault_upcall;
38      env_run(curenv);
39
40  }
41
42  // Destroy the environment that caused the fault.
43  cprintf("[%08x] user fault va %08x ip %08x\n",
44      curenv->env_id, fault_va, tf->tf_eip);
45  print_trapframe(tf);
46  env_destroy(curenv);
47  }

```

## Exercise 10

`_pgfault_upcall` 的如下代码实现了跳回 `env` 发生 page fault 的地方，继续运行。

```

1  movl 0x30(%esp), %ecx
2  subl $4, %ecx
3  movl %ecx, 0x30(%esp)
4  movl 0x28(%esp), %edx
5  movl %edx, (%ecx)

```



---

```
6
7     addl $8, %esp
8     popal
9
10    addl $4, %esp
11    popfl
12
13    pop %esp
14
15    ret
16
```

## Exercise 11

set\_pgfault\_handler 设置了 fault handler.

```
1  void
2  set_pgfault_handler(void (*handler)(struct UTrapframe *utf))
3  {
4      int r;
5
6      if (_pgfault_handler == 0) {
7          // First time through!
8          // LAB 4: Your code here.
9          int e = sys_page_alloc(thisenv->env_id, (void *) (UXSTACKTOP - PGSIZE), PTE_U|PTE_P|PTE_W);
10         if (e<0)
11             panic("set_pgfault_handler failed!!!");
12         e = sys_env_set_pgfault_upcall(thisenv->env_id, _pgfault_upcall);
13         if (e<0)
14             panic("set_pgfault_handler failed!!!");
15
16     }
17
18     // Save handler pointer for assembly to call.
19     _pgfault_handler = handler;
20 }
21
```

## Exercise 12

fork 实现了一个 copy-on-write 的 fork.

---

```

1  envid_t
2  fork(void)
3  {
4      int r;
5      envid_t id = sys_getenvid();
6      set_pgfault_handler(pgfault);
7      envid_t t = sys_exofork();
8      if (t<0)
9          panic("fork: exofork error!!!");
10     if (t == 0){
11         thisenv = (&envs[ENVX(sys_getenvid())]);
12     }
13     else{
14         for (uint32_t addr = 0;addr<USTACKTOP; addr+=PGSIZE)
15             if ((uvpd[PDX(addr)] & PTE_P) == PTE_P && (uvpt[PGNUM(addr)] & PTE_P) == PTE_P)
16                 duppage(t, addr/PGSIZE);
17         r = sys_page_alloc(t, (void *) (UTOP-PGSIZE), PTE_U | PTE_P | PTE_W);
18         if (r<0)
19             panic("fork: %e", r);
20         void _pgfault_upcall();
21         r = sys_env_set_pgfault_upcall(t, _pgfault_upcall);
22         if (r<0)
23             panic("set_pgfault_handler failed!!!");
24         r = sys_env_set_status(t, ENV_RUNNABLE);
25         if (r<0)
26             panic("fork: %e", r);
27     }
28     return t;
29 }

```

duppage 把当前 env 的第 pn 个虚拟页 map 到另一个 env 的相同虚拟位置. 如果当前页是 copy on write 或 writable, 那么将页映射都设为 copy on write.

```

1  static int
2  duppage(envid_t envid, unsigned pn)
3  {
4      int r;
5      void * va = (void *) (pn * PGSIZE);
6      envid_t id = sys_getenvid();
7      if ((uvpt[pn] & PTE_W) == PTE_W || (uvpt[pn] & PTE_COW) == PTE_COW){
8          r = sys_page_map(id, va, envid, va, PTE_COW | PTE_U | PTE_P);

```

---

```

9     if (r<0)
10         return r;
11     r = sys_page_map(id, va, id, va, PTE_COW | PTE_U | PTE_P);
12     if (r<0)
13         return r;
14 }
15 else{
16     r = sys_page_map(id, va, envid, va, PTE_U | PTE_P);
17     if (r<0)
18         return r;
19 }
20 return 0;
21 }

```

这个是真正的 page fault handler, 处理了 copy on write 的情况.

```

1  static void
2  pgfault(struct UTrapframe *utf)
3  {
4      void *addr = (void *) utf->utf_fault_va;
5      uint32_t err = utf->utf_err;
6      int r;
7
8      if ( ((uvpt[(uint32_t)addr/PGSIZE] & PTE_COW) == PTE_COW) && ((err & FEC_WR) == FEC_WR) ){
9
10     }
11     else
12         panic("pgfault: not a write to a copy-on-write page!!!");
13
14     envid_t id = sys_getenvid();
15     r = sys_page_alloc(id, PFTEMP, PTE_U | PTE_W | PTE_P);
16     if (r<0)
17         panic("pgfault: %e", r);
18     memcpy(PFTEMP, ROUNDDOWN(addr, PGSIZE), PGSIZE);
19     r = sys_page_map(id, PFTEMP, id, ROUNDDOWN(addr, PGSIZE), PTE_U | PTE_W | PTE_P);
20     if (r<0)
21         panic("pgfault: %e", r);
22     r = sys_page_unmap(id, PFTEMP);
23     if (r<0)
24         panic("pgfault: %e", r);
25 }

```

---

## Exercise 13

就和 lab3 的时候设立 IDT 的方法一样，现在在上面又加了一些。

```
1 TRAPHANDLER_NOEC(TIMER, IRQ_OFFSET + IRQ_TIMER)
2 TRAPHANDLER_NOEC(KBD, IRQ_OFFSET + IRQ_KBD)
3 TRAPHANDLER_NOEC(SERIAL, IRQ_OFFSET + IRQ_SERIAL)
4 TRAPHANDLER_NOEC(SPURIOUS, IRQ_OFFSET + IRQ_SPURIOUS)
5 TRAPHANDLER_NOEC(IDE, IRQ_OFFSET + IRQ_IDE)
6 TRAPHANDLER_NOEC(ERROR, IRQ_OFFSET + IRQ_ERROR)

1 void TIMER();
2 SETGATE(idt[IRQ_OFFSET + IRQ_TIMER], 0, GD_KT, TIMER, 0);
3 void KBD();
4 SETGATE(idt[IRQ_OFFSET + IRQ_KBD], 0, GD_KT, KBD, 0);
5 void SERIAL();
6 SETGATE(idt[IRQ_OFFSET + IRQ_SERIAL], 0, GD_KT, SERIAL, 0);
7 void SPURIOUS();
8 SETGATE(idt[IRQ_OFFSET + IRQ_SPURIOUS], 0, GD_KT, SPURIOUS, 0);
9 void IDE();
10 SETGATE(idt[IRQ_OFFSET + IRQ_IDE], 0, GD_KT, IDE, 0);
11 void ERROR();
12 SETGATE(idt[IRQ_OFFSET + IRQ_ERROR], 0, GD_KT, ERROR, 0);
```

## Exercise 14

在 trap\_dispatch 里增加了对 timer interrupt 的处理。

```
1 if (tf->tf_trapno == IRQ_OFFSET + IRQ_TIMER) {
2     lapic_eoi();
3     sched_yield();
4     return;
5 }
```

## Exercise 15

在 sys\_ipc\_recv 这个 syscall 将 env 设为在 recv 的状态, 并将 env\_status 设为 ENV\_NOT\_RUNNABLE 使得在收到东西前不去跑这个 env。

```
1 static int
2 sys_ipc_recv(void *dstva)
3 {
```

---

```

4  if ((uint32_t)dstva < UTOP && (uint32_t)dstva%PGSIZE!=0 )
5      return -E_INVALID;
6  curenv->env_ipc_recving = 1;
7  curenv->env_ipc_dstva = dstva;
8  curenv->env_status = ENV_NOT_RUNNABLE;
9  return 0;
10 }
```

在 sys\_ipc\_try\_send 这个 syscall 尝试向一个 env 发送东西.

```

1  static int
2  sys_ipc_try_send(envid_t envid, uint32_t value, void *srcva, unsigned perm)
3  {
4      // LAB 4: Your code here.
5      int r;
6      struct Env * env;
7      struct PageInfo * pp;
8      pte_t * pte;
9      if ((r=envid2env(envid,&env,0))<0)
10         return r;
11
12     if (env->env_ipc_recving != 1)
13         return -E_IPC_NOT_RECV;
14     if ((uint32_t)srcva < UTOP && (uint32_t)srcva%PGSIZE!=0)
15         return -E_INVALID;
16
17     if ((uint32_t)srcva < UTOP && ((perm & (PTE_U | PTE_P)) != (PTE_U | PTE_P) || (((perm ^ (PTE_U | PTE_P)) & 0x3) != 0)))
18         return -E_INVALID;
19     if ((uint32_t)srcva < UTOP && (pp = page_lookup(curenv->env_pgdir, srcva, &pte)) == NULL)
20         return -E_INVALID;
21
22     if ((uint32_t)srcva < UTOP && ((*pte)&PTE_W) == 0 && (perm & PTE_W)!=0)
23         return -E_INVALID;
24
25     if ((uint32_t)(env->env_ipc_dstva) < UTOP && (uint32_t)srcva < UTOP){
26         if ((r = page_insert(env->env_pgdir, pp, env->env_ipc_dstva, perm))<0)
27             return r;
28         env->env_ipc_perm = perm;
29     }
30     else
31         env->env_ipc_perm = 0;
```

---

```

32
33     env->env_ipc_recving = 0;
34     env->env_ipc_from = curenv->env_id;
35     env->env_ipc_value = value;
36     env->env_status = ENV_RUNNABLE;
37     return 0;
38 }

```

在 ipc\_recv 用来收东西，基本直接调用 sys\_ipc\_recv.

```

1  int32_t
2  ipc_recv(envid_t *from_env_store, void *pg, int *perm_store)
3  {
4      if (pg == NULL)
5          pg = (void *)UTOP;
6      int r;
7      if ((r = sys_ipc_recv(pg)) < 0){
8          if (from_env_store) (*from_env_store) = 0;
9          if (perm_store) (*perm_store)=0;
10         return r;
11     }
12     if (from_env_store) (*from_env_store) = thisenv->env_ipc_from;
13     if (perm_store) (*perm_store) = thisenv->env_ipc_perm;
14     return thisenv->env_ipc_value;
15 }

```

在 ipc\_send 用来发东西,基本直接调用 sys\_ipc\_send, 在对面不处在收东西的状态的时候就先 sys\_yield 调度别的进程，减少了对 CPU 资源的浪费.

```

1  void
2  ipc_send(envid_t to_env, uint32_t val, void *pg, int perm)
3  {
4      int r;
5      if (pg == NULL)
6          pg = (void *)UTOP;
7      while (1){
8          r=sys_ipc_try_send(to_env, val, pg, perm);
9          if (r<0 && r!=-E_IPC_NOT_RECV)
10             panic("ipc_send: %e", r);
11         if (r==0) break;
12         sys_yield();
13     }
14 }

```