

Report [Lab 4: Preemptive Multitasking]

I complete the first challenge after exercise 15, completing an `ipc_send()` with no loops.

Part A: Multiprocessor Support and Cooperative Multitasking

Multiprocessor Support

Exercise 1

`mmio_map_region()` is used in `lapic_init()` to map LAPIC's 4K MMIO region into virtual address space for the convenience of accessing. Because back the pages between MMIOBASE and MMIOLIM are no normal page frames, their content can't be cached and must apply write-through policy. This is done by setting specify permission bits in the PTE by `boot_map_region()`

```
1 size = ROUNDUP(size, PGSIZE);
2 if (base + size > MMIOLIM)
3 {
4     panic("mmio_map_region: exceed MMIOLIM");
5 }
6 boot_map_region(kern_pgdir, base, size, pa, PTE_W | PTE_PCD | PTE_PWT);
7 void * re = (void *)base;
8 base += size;
9 return re;
```

Application Processor Bootstrap

Exercise 2

`boot_aps()` is called from `i386_init()` to awake all application processors.

First it copies the assembly code in `kern/mpentry.S` to the physical address `MPENTRY_PADDR` (a low address to allow the code to run in real mode). Then each AP is given their own kernel stack to execute the assembly code starting at `mpentry_start` by calling `lapic_statap()`. Then it spins to wait for the AP startup before awaking the next.

In `mpentry.S`, things are quite similar to the BSP bootstrap. First protected mode is turned on then GDT is loaded and paging is enabled with `entry_pgdir`. Finally the control flow jumps to `mp_main()`

In `mp_main`, `kernel_pgdir` is loaded and then `lapic_init()`, `env_init_percpu()`, `trap_init_percpu()` is called to initialize local config. Then it set `thiscpu->cpu_status` as `CPU_STARTED` to announce this.

So we change `page_init()` in `pmap.c` to reserve space for AP start code.

```
1 if (i == PGNUM(MPENTRY_PADDR)) // reserved for AP start code
2     continue;
```

Question 1

`MPBOOTPHYS` is intended to interpret the address of the symbols `mpentry.S` uses e.g. `gdt` from kernel virtual address to the physical address at which the `mpentry.S` actually runs i.e. above `0x7000`

`kern/mpentry.S` and `boot/boot.S` are quite the same in their content but they are linked to different address by the linker. We can observe this using `objdump`, i.e they've got different VMA

```
1 $ objdump -x obj/kern/kernel
2 # ...
3 Sections:
4 Idx Name          Size      VMA      LMA      File off  Algn
5 0 .text           00005d21 f0100000 00100000 00001000 2**4
6                      CONTENTS, ALLOC, LOAD, READONLY, CODE
7 # ...
8 $ objdump -x obj/boot/boot.o
9 # ...
10 Sections:
11 Idx Name          Size      VMA      LMA      File off  Algn
12 0 .text           0000006a 00000000 00000000 00000034 2**2
13                      CONTENTS, ALLOC, LOAD, RELOC, READONLY, CODE
14 # ...
```

So if we omit using `MPBOOTPHYS` to translate symbols in `kern/mpentry.S`, the code will try to address those symbols using high virtual addresses when we are still in real mode. That's definitely what we expected.

```
1 $ objdump -x obj/kern/kernel | grep gdt
2 f0105360 l          .text 00000000 gdt desc
3 f0105348 l          .text 00000000 gdt
4 f0124340 g          0 .data 00000068 gdt
5 f0124320 g          0 .data 00000006 gdt_pd
```

Per-CPU State and Initialization

Exercise 3

In the past we only map `bootstack` to back-store BSP's kernel stack. Now we map `percpu_kstacks` to back-store each CPU's kernel stack because they share the same address space while needing to handle interrupts simultaneously.

```
1     for (int i = 0; i < NCPU; i++)
2     {
3         intptr_t kstacktop_i = KSTACKTOP - i * (KSTKSIZE + KSTKGAP);
4         boot_map_region(kern_pgdir, kstacktop_i - KSTKSIZE, KSTKSIZE,
5                         PADDR(percpu_kstacks[i]), PTE_W);
6     }
```

Exercise 4

For each CPU, a TSS and a TSS descriptor is needed to indicate the stack position for stack switch happening at interrupt handling which changes CPL. We set `esp0` and `ss0` field as the kernel stack we preserved for `thiscpu` and modify the entry in GDT to point to the `cpu_ts` structure we initialized.

There is one field in TSS called `ts_iomb`, which points to the beginning of the I/O permission bit map and the end of the interrupt redirection bit map, that are stored in TSS at higher addresses if present. [1] We set it to end of TSS basic structure (which indicates the non-existence of interrupt redirection bit map) with the `limit` field in GDT entry set as `sizeof(struct Taskstate) - 1` (which indicates the non-existence of I/O permission bit map). If `ts_iomb` is set to 0, some range of address might be misinterpreted as I/O permission bit map, which is dangerous.

Then we load the TSS descriptor using `ltr` and the same IDT descriptor as BSP using `lidt`.

```
1  int i = cpunum();
2  intptr_t kstacktop_i = KSTACKTOP - i * (KSTKSIZE + KSTKGAP);
3  thiscpu->cpu_ts.ts_esp0 = kstacktop_i;
4  thiscpu->cpu_ts.ts_ss0 = GD_KD;
5  thiscpu->cpu_ts.ts_iomb = sizeof(struct Taskstate);
6
7  gdt[(GD_TSS0 >> 3) + i] = SEG16(STS_T32A, (uint32_t) (&thiscpu->cpu_ts),
8                                sizeof(struct Taskstate) - 1, 0);
9  gdt[(GD_TSS0 >> 3) + i].sd_s = 0;
10
11  ltr(GD_TSS0 + (i << 3));
12
13  lidt(&idt_pd);
14
15  return;
```

Locking

Exercise 5

Apply the big kernel lock at the indicated positions

```
1  // i386_init()
2  // Acquire the big kernel lock before waking up APs
3  // Your code here:
4  lock_kernel();
```

```
1  // mp_main()
2  // Now that we have finished some basic setup, call sched_yield()
3  // to start running processes on this CPU. But make sure that
4  // only one CPU can enter the scheduler at a time!
5  //
6  // Your code here:
7  lock_kernel();
8  sched_yield();
```

```

1 // trap()
2     // Trapped from user mode.
3     // Acquire the big kernel lock before doing any
4     // serious kernel work.
5     // LAB 4: Your code here.
6     assert(curenv);
7     lock_kernel();

```

```

1 // env_run()
2     if (curenv != NULL)
3     {
4         if (curenv->env_status == ENV_RUNNING)
5         {
6             curenv->env_status = ENV_RUNNABLE;
7         }
8         // we have saved context of curenv at _alltrap
9     }
10    curenv = e;
11    e->env_status = ENV_RUNNING;
12    e->env_runs++;
13    lcr3(PADDR(e->env_pgdir));
14
15    unlock_kernel();
16
17    env_pop_tf(&e->env_tf);

```

Question 2

Each CPU needs its own stack to store the information because it might be blocked during execution in kernel mode. For example if the kernel code in one CPU performs some IO then it's blocked, some other CPU will acquire the big kernel lock and run its kernel code, corrupting the content of the shared stack.

Round-Robin Scheduling

Exercise 6

First we acquire the index of current running environment by `ENVX(curenv->envid)`, or beginning from 0 if `curenv` is `NULL`. If there is no other environment `RUNNABLE` and `curenv` is `ENV_RUNNING`, we keep it as `curenv`.

```

1 // LAB 4: Your code here.
2 for (int i = 0, j = curenv == NULL ? 0 : ENVX(curenv->env_id); i < NENV;
   i++, j++)
3 {
4     if (j >= NENV) j -= NENV;
5     if (envs[j].env_status == ENV_RUNNABLE)
6     {
7         env_run(&envs[j]);
8     }
9 }
10 if (curenv != NULL && curenv->env_status == ENV_RUNNING)
11 {
12     env_run(curenv);
13 }

```

Dispatch `sys_yield()` in `syscall()`. Note that it never returns.

```

1 case SYS_yield:
2     sys_yield();

```

create 3 environments running `user/yield.c` to test.

```

1 // Touch all you want.
2 // ENV_CREATE(user_primes, ENV_TYPE_USER);
3 ENV_CREATE(user_yield, ENV_TYPE_USER);
4 ENV_CREATE(user_yield, ENV_TYPE_USER);
5 ENV_CREATE(user_yield, ENV_TYPE_USER);

```

```

1 hiesa@IL:~/workspace/6.828/lab$ make qemu-nox CPUS=2
2 # ...
3 [00000000] new env 00001000
4 [00000000] new env 00001001
5 [00000000] new env 00001002
6 Hello, I am environment 00001000.
7 Hello, I am environment 00001001.
8 Back in environment 00001000, iteration 0.
9 Hello, I am environment 00001002.
10 Back in environment 00001001, iteration 0.
11 Back in environment 00001000, iteration 1.
12 Back in environment 00001002, iteration 0.
13 Back in environment 00001001, iteration 1.
14 Back in environment 00001000, iteration 2.
15 Back in environment 00001002, iteration 1.
16 # ...

```

Question 3

The pointer `e` points one element in the kernel data structure `envs`. Since all environments map the kernel part of virtual address space in the same manner, `e` points to the identical data before and after the addressing switch.

Question 4

Because environments are oblivious to context switch, if the old registers are not stored and later recovered, environment will run with corrupted registers, giving out unexpected result.

This is down by two parts. First part is down by x86 hardware, pushing program counter and stack register into the stack; second part is down at `_alltrap` and `trap()`, where all general registers are pushed onto stack forming a `Trapframe` structure and copied into `curenv->env_tf`.

System Calls for Environment Creation

Exercise 7

In this exercise we will implement a number of system calls for user environments.

```
1  case SYS_exofork:
2      return sys_exofork();
3  case SYS_env_set_status:
4      return sys_env_set_status(a1, a2);
5  case SYS_page_alloc:
6      return sys_page_alloc(a1, (void *)a2, a3);
7  case SYS_page_map:
8      return sys_page_map(a1, (void *)a2, a3, (void *)a4, a5);
9  case SYS_page_unmap:
10     return sys_page_unmap(a1, (void *)a2);
```

SYS_exofork

This function is to create a new environment which is basically what `env_alloc` creates except that its status is set as `ENV_NOT_RUNNABLE` for setups and its register values are copied from the parent environment.

We return the `env_id` of child environment while tweaking the the value of `%eax` of the child environment to 0 so the two environment involved see different return value.

```
1  // LAB 4: Your code here.
2  int retval;
3  struct Env *child;
4
5  if ((retval = env_alloc(&child, curenv->env_id)))
6  {
7      return retval;
8  }
9
10 child->env_status = ENV_NOT_RUNNABLE;
11 child->env_tf = curenv->env_tf;
12 child->env_tf.tf_regs.reg_eax = 0;
13
14 return child->env_id;
```

SYS_env_set_status

The main effort in this and following functions are put into checking the validity of the arguments.

First we have to ensure the environment which calls this function i.e. `curenv` has the permission to manipulate the context of the environment indicated by `envid`. This is done in `envid2env()`, as whether the environment indicated by `envid` (if still existing) is `curenv` itself or its direct child.

We also checks whether the status is not `ENV_RUNNABLE` or `ENV_NOT_RUNNABLE`.

```
1   if (status != ENV_RUNNABLE && status != ENV_NOT_RUNNABLE)
2   {
3       return -E_INVAL;
4   }
5
6   struct Env* e;
7   int retval;
8
9   if ((retval = envid2env(envid, &e, true)))
10  {
11      return retval;
12  }
13
14  e->env_status = status;
15
16  return 0;
```

SYS_page_alloc

In this function we allocate a physical page with `page_alloc()` and insert it at the address indicated by `va` using `page_insert()` with `perm` as the permission, after we have checked that the `va` is under `UTOP` and page_aligned and `perm` is valid (present bit and user bit are set while bits other than those set in `PTE_SYSCALL` aren't set).

```
1   // LAB 4: Your code here.
2   if ((intptr_t)(va) >= UTOP || (intptr_t)(va) % PGSIZE)
3   {
4       return -E_INVAL;
5   }
6   if ((~perm & PTE_P) || (~perm & PTE_U) || (perm & ~PTE_SYSCALL))
7   {
8       return -E_INVAL;
9   }
10
11  struct Env *e;
12  struct PageInfo *p;
13  int retval;
14
15  if ((retval = envid2env(envid, &e, true)))
16  {
17      return retval;
18  }
19
```

```

20     p = page_alloc(ALLOC_ZERO);
21     if (p == NULL)
22     {
23         return -E_NO_MEM;
24     }
25     if ((retval = page_insert(e->env_pgdir, p, va, perm)))
26     {
27         page_free(p);
28         return -E_NO_MEM;
29     }
30
31     return 0;

```

SYS_page_map

In this function we map the physical page indicated by `srcva` in source environment's address space (which can be queried by `page_lookup()`) then map it at `dstva` in destination environment using `page_insert` with new permission bits `perm`.

Same checks like those above are performed upon both source and destination `va` and `envid`. In addition, we have to ensure that a read-only page in the source can't be mapped at the destination with the write bit set. This is done by checking the PTE in source virtual address space which the `page_lookup()` stored

```

1 // LAB 4: Your code here.
2 struct Env *srce, *dste;
3 struct PageInfo *p;
4 pte_t *pte;
5 int retval;
6
7 if ((retval = envid2env(srcenvid, &srce, true)))
8 {
9     return retval;
10 }
11 if ((retval = envid2env(dstenvid, &dste, true)))
12 {
13     return retval;
14 }
15
16 if ((intptr_t)(srcva) >= UTOP || (intptr_t)(srcva) % PGSIZE)
17 {
18     return -E_INVALID;
19 }
20 if ((intptr_t)(dstva) >= UTOP || (intptr_t)(dstva) % PGSIZE)
21 {
22     return -E_INVALID;
23 }
24
25 p = page_lookup(srce->env_pgdir, srcva, &pte);
26 if (p == NULL)
27 {
28     return -E_INVALID;
29 }
30 if ((~perm & PTE_P) || (~perm & PTE_U) || (perm & ~PTE_SYSCALL))
31 {

```



```

32         return -E_INVALID;
33     }
34     if ((perm & PTE_W) && (~(*pte) & PTE_W))
35     {
36         return -E_INVALID;
37     }
38
39     if ((retval = page_insert(dste->env_pgdir, p, dstva, perm)))
40     {
41         return retval;
42     }
43
44     return 0;

```

SYS_page_unmap

In this function, we check the validity of the arguments and then use `page_remove()` to remove the physical page at `va` in the address space of the environment with `envid`.

If there wasn't any page at `va`, the function silently succeed, like the behavior of `page_remove()` it calls.

```

1 // LAB 4: Your code here.
2 struct Env *e;
3 int retval;
4
5 if ((retval = envid2env(envid, &e, true)))
6 {
7     return retval;
8 }
9 if (((intptr_t)(va) >= UTOP || (intptr_t)(va) % PGSIZE)
10 {
11     return -E_INVALID;
12 }
13
14 page_remove(e->env_pgdir, va);
15
16 return 0;

```

Part B: Copy-on-Write Fork

User-level page fault handling

Setting the Page Fault Handler

Exercise 8

We check the validity of `envid` before we set up `env_pgfault_upcall` like what we do in the previous exercise.

```

1 static int
2 sys_env_set_pgfault_upcall(envid_t envid, void *func)
3 {
4     // LAB 4: Your code here.

```

```

5     struct Env *e;
6     int retval;
7
8     if ((retval = envid2env(envid, &e, true)))
9     {
10         return retval;
11     }
12
13     e->env_pgfault_upcall = func;
14
15     return 0;
16
17     // panic("sys_env_set_pgfault_upcall not implemented");
18 }

```

```

1     case SYS_env_set_pgfault_upcall:
2         return sys_env_set_pgfault_upcall(a1, (void *)a2);

```

Invoking the User Page Fault Handler

Exercise 9

If current environment has set up its `env_pgfault_upcall`, we should transfer control to user's registered handler when page fault occurs. To provide information for user defined page fault handler to recover the fault and return to previous instruction causing the fault, we need to push a `UTrapframe` structure onto user exception stack.

There are two cases which can be determined by the value of `%esp` at trap time:

1. The page fault occurs when user environment is running with normal stack. Then we should switch the stack to start just below `UXSTACKTOP`
2. The page fault occurs during the handling of another fault, where we should switch to a new stack below current exception stack position i.e. trap-time `%esp`, leaving a 4-word gap for scratch space. (used by the assembly language stub to help page fault handle return)

Either case, we should make sure there is enough available space for us to store the `UTrapframe` structure

```

1     if (curenv->env_pgfault_upcall)
2     {
3         struct UTrapframe *utf;
4
5         if (tf->tf_esp >= UXSTACKTOP - PGSIZE && tf->tf_esp < UXSTACKTOP)
6         {
7             // already running on the user exception stack
8             utf = (struct UTrapframe*)(tf->tf_esp - 4 - sizeof(struct
UTrapframe));
9         }
10        else
11        {
12            // running on normal user stack
13            utf = (struct UTrapframe*)(UXSTACKTOP - sizeof(struct
UTrapframe));

```

```

14     }
15
16     // validity check of exception stack
17     user_mem_assert(curenv, utf, sizeof(struct UTrapframe), PTE_W);

```

Then we fill in the fields of `UTrapFrame` and change `tf_eip` and `tf_esp` so that the environment will run the registered page fault handler with a new stack by `env_run()`

```

1     // fill in utf
2     utf->utf_fault_va = fault_va;
3     utf->utf_err = tf->tf_err;
4     utf->utf_regs = tf->tf_regs;
5     utf->utf_eip = tf->tf_eip;
6     utf->utf_eflags = tf->tf_eflags;
7     utf->utf_esp = tf->tf_esp;
8
9     // run the user page fault handler with new stack
10    tf->tf_eip = (intptr_t)curenv->env_pgfault_upcall;
11    tf->tf_esp = (uintptr_t)utf;
12    env_run(curenv);

```

User-mode Page Fault Entrypoint

Exercise 10

In this exercise we will complete the assembly to return to the instruction which invokes the page fault from page fault handler, which means we have to restore the execution state with much precaution.

First we need to restore the general registers, after which we can no longer use them to store temporary values. Then we need to restore `eflags` using `popfl`, after which we can no longer use arithmetic operations. Finally we have to restore `%esp` and `%eip`, which is the most sophisticated phase.

Since all general registers are restored, we can't use `jmp` since it needs a register to store the address. Neither can we simply use `ret` on exception stack otherwise the process will run with a wrong `%esp`.

So we will do some tricks, pushing the trap time `%eip` we want to return to onto the *trap time* stack and let `%esp` point to it before we use `ret`. That will set `%esp` and `%eip` both in correct place.

```

1     // LAB 4: Your code here.
2     movl 40(%esp), %eax // trap-time eip
3     subl $4, 48(%esp)  // tweak the trap-time stack for later ret
4     movl 48(%esp), %edx
5     movl %eax, (%edx)  // use 4 bytes below trap-time stack as scratch space
    for later ret
6
7     // Restore the trap-time registers. After you do this, you
8     // can no longer modify any general-purpose registers.
9     // LAB 4: Your code here.
10    addl $8, %esp // now point to utf_regs
11    popal

```

```

12
13 // Restore eflags from the stack. After you do this, you can
14 // no longer use arithmetic operations or anything else that
15 // modifies eflags.
16 // LAB 4: Your code here.
17 addl $4, %esp // now point to utf_eflags
18 popfl
19
20 // Switch back to the adjusted trap-time stack.
21 // LAB 4: Your code here.
22 movl (%esp), %esp
23
24 // Return to re-execute the instruction that faulted.
25 // LAB 4: Your code here.
26 ret

```

Exercise 11

In the user library function `set_pgfault_handler`, we use a global variable to store the user provided page fault handler. If it's called the first time, it will allocate a page for user exception stack and register the assembly page fault handler entry point with a system call.

```

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         sys_page_alloc(0, (void *) (UXSTACKTOP - PGSIZE), PTE_W | PTE_U |
PTE_P);
10        sys_env_set_pgfault_upcall(0, _pgfault_upcall);
11        // panic("set_pgfault_handler not implemented");
12    }
13
14    // Save handler pointer for assembly to call.
15    _pgfault_handler = handler;
16 }

```

Implementing Copy-on-Write Fork

Exercise 12

In this exercise we will exploit the user-level page fault handler we implemented above to complete the function of copy-on-write fork.

The 11th bit of PTE is left available for software so we use it as the COW bit. We implement a user level page fault handler to deal with page fault which is a write access to a COW page.

```

1 // LAB 4: Your code here.
2 pte_t *pte = PGADDR(PDX(UVPT), PDX(addr), PTX(addr) << 2);
3
4 if (~err & FEC_WR)
5 {
6     panic("pgfault: not write access");
7 }
8 if (~(*pte) & PTE_COW)
9 {
10    panic("pgfault: not access to copy-on-write page");
11 }
12

```

We use the page at `PFTEMP` which is below `UTEXT` as the scratch space. We allocate a temporary page there and copy the content of the fault page to it. Then we unmap the COW page and map the new private page frame at the fault address.

```

1 // LAB 4: Your code here.
2 if ((r = sys_page_alloc(0, (void *)PFTEMP, PTE_P | PTE_U | PTE_W)))
3 {
4     panic("pgfault: %e\n", r);
5 }
6 memcpy((void *)PFTEMP, ROUNDDOWN(addr, PGSIZE), PGSIZE);
7 if ((r = sys_page_unmap(0, ROUNDDOWN(addr, PGSIZE))))
8 {
9     panic("pgfault: %e\n", r);
10 }
11 if ((r = sys_page_map(0, (void *)PFTEMP, 0, ROUNDDOWN(addr, PGSIZE),
12 PTE_P | PTE_U | PTE_W)))
13 {
14     panic("pgfault: %e\n", r);
15 }
16

```

In `fork()`, we register the COW page fault handler first and use `sys_exofork()` to create the child environment with a virtual address which is current empty below `UTOP`. Then we use `duppage()` to copy all pages existing in current environment under `UTOP` to the child environment expect of the page of user exception stack because if it's marked as copy-on-write, no one would do the work for it.

After that we register the same page fault handler for the child (as now the code of the page fault handler also resides at the same address in child's virtual address space) and allocate a page for its user exception stack much like what we do when we first enter `set_pgfault_handler()`. Then things are done and we mark the child as `ENV_RUNNABLE`

```

1 envid_t
2 fork(void)
3 {
4     // LAB 4: Your code here.
5     envid_t child;
6     int r;
7
8     set_pgfault_handler(&pgfault);
9

```

```

10     if ((child = sys_exofork()) < 0)
11     {
12         panic("fork: %e", child);
13     }
14
15     // I'm the child
16     if (child == 0)
17     {
18         thisenv = &envs[ENVX(sys_getenvid())];
19         return child;
20     }
21
22     for (int i = 0; i <= PDX(USTACKTOP); i++)
23     {
24         pde_t *pde = PGADDR(PDX(UVPT), PDX(UVPT), (i << 2));
25         if (~(*pde) & PTE_P) continue;
26         for (int j = 0; j < (PGSIZE >> 2); j++)
27         {
28             if (i == PDX(USTACKTOP) && j >= PTX(USTACKTOP))
29             {
30                 break;
31             }
32             pte_t *pte = PGADDR(PDX(UVPT), i, j << 2);
33             if (~(*pte) & PTE_P) continue;
34
35             duppage(child, i * (PGSIZE >> 2) + j);
36         }
37     }
38
39     extern void _pgfault_upcall(void);
40     if ((r = sys_page_alloc(child, (void *) (UXSTACKTOP - PGSIZE), PTE_W |
PTE_U | PTE_P)))
41     {
42         panic("fork: %e\n", r);
43     }
44     if ((r = sys_env_set_pgfault_upcall(child, _pgfault_upcall)))
45     {
46         panic("fork: %e\n", r);
47     }
48
49     if ((r = sys_env_set_status(child, ENV_RUNNABLE)))
50     {
51         panic("fork: %e\n", r);
52     }
53
54     return child;
55
56     // panic("fork not implemented");
57 }
58

```

In `duppage()`, we simply maps the read-only pages with the same `perm` to the other environment while for writable or copy-on-write pages, we map the page in the new environment as a COW page and then mark it as COW in current address space.

it's important to mark a page as COW in the child before marking it in the parent because if we mark the page in parent first, there would be a race between writing to the page in the parent and marking it COW in the child. For example if the page stands for the current stack we might invoke the page fault handler when pushing the arguments onto the stack to call the `sys_page_map` so that the page is turned back to a normal writable page before it's marked COW in the child. Then the same page frame will be writable in the parent while COW in the child so the write in the parent would affect what the child will read. What a disaster!

For the same reason we should always mark the page COW in the parent even it's COW at the beginning, as it might turn to a writable page in the procedure.

```
1 static int
2 duppage(envid_t envid, unsigned pn)
3 {
4     int r;
5
6     // LAB 4: Your code here.
7     pte_t *pte = PGADDR(PDX(UVPT), pn / (PGSIZE >> 2), (pn % (PGSIZE >> 2))
8 << 2);
9     void *addr = (void *) (pn * PGSIZE);
10
11     if ((*pte) & (PTE_W | PTE_COW))
12     {
13         // writable or copy-on-write page
14         int perm = ((*pte) & PTE_SYSCALL) & (~PTE_W) | PTE_COW;
15         if ((r = sys_page_map(0, addr, envid, addr, perm)))
16         {
17             panic("duppage: %e\n", r);
18         }
19         if ((r = sys_page_map(0, addr, 0, addr, perm)))
20         {
21             panic("duppage: %e\n", r);
22         }
23     }
24     else
25     {
26         // read-only page
27         int perm = (*pte) & PTE_SYSCALL;
28         if ((r = sys_page_map(0, addr, envid, addr, perm)))
29         {
30             panic("duppage: %e\n", r);
31         }
32     }
33     // panic("duppage not implemented");
34     return 0;
35 }
```

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

Clock Interrupts and Preemption

Interrupt discipline

Exercise 13

We do the same thing as in previous lab.

```
1 // trapentry.S
2 TRAPHANDLER_NOEC(ENTRY_IRQ0, IRQ_OFFSET + 0)
3 TRAPHANDLER_NOEC(ENTRY_IRQ1, IRQ_OFFSET + 1)
4 TRAPHANDLER_NOEC(ENTRY_IRQ2, IRQ_OFFSET + 2)
5 TRAPHANDLER_NOEC(ENTRY_IRQ3, IRQ_OFFSET + 3)
6 TRAPHANDLER_NOEC(ENTRY_IRQ4, IRQ_OFFSET + 4)
7 TRAPHANDLER_NOEC(ENTRY_IRQ5, IRQ_OFFSET + 5)
8 TRAPHANDLER_NOEC(ENTRY_IRQ6, IRQ_OFFSET + 6)
9 TRAPHANDLER_NOEC(ENTRY_IRQ7, IRQ_OFFSET + 7)
10 TRAPHANDLER_NOEC(ENTRY_IRQ8, IRQ_OFFSET + 8)
11 TRAPHANDLER_NOEC(ENTRY_IRQ9, IRQ_OFFSET + 9)
12 TRAPHANDLER_NOEC(ENTRY_IRQ10, IRQ_OFFSET + 10)
13 TRAPHANDLER_NOEC(ENTRY_IRQ11, IRQ_OFFSET + 11)
14 TRAPHANDLER_NOEC(ENTRY_IRQ12, IRQ_OFFSET + 12)
15 TRAPHANDLER_NOEC(ENTRY_IRQ13, IRQ_OFFSET + 13)
16 TRAPHANDLER_NOEC(ENTRY_IRQ14, IRQ_OFFSET + 14)
17 TRAPHANDLER_NOEC(ENTRY_IRQ15, IRQ_OFFSET + 15)
```

```
1 // trap.c
2 // entry point for IRQ 0 ~ 15
3 extern void ENTRY_IRQ0();
4 extern void ENTRY_IRQ1();
5 extern void ENTRY_IRQ2();
6 extern void ENTRY_IRQ3();
7 extern void ENTRY_IRQ4();
8 extern void ENTRY_IRQ5();
9 extern void ENTRY_IRQ6();
10 extern void ENTRY_IRQ7();
11 extern void ENTRY_IRQ8();
12 extern void ENTRY_IRQ9();
13 extern void ENTRY_IRQ10();
14 extern void ENTRY_IRQ11();
15 extern void ENTRY_IRQ12();
16 extern void ENTRY_IRQ13();
17 extern void ENTRY_IRQ14();
18 extern void ENTRY_IRQ15();
```

```
1 // trap_init() in trap.c
2 // set up IDT entries for IRQ 0 ~ 15
3 SETGATE(idt[IRQ_OFFSET + 0], 0, GD_KT, &ENTRY_IRQ0, 0);
4 SETGATE(idt[IRQ_OFFSET + 1], 0, GD_KT, &ENTRY_IRQ1, 0);
5 SETGATE(idt[IRQ_OFFSET + 2], 0, GD_KT, &ENTRY_IRQ2, 0);
```



```

6   SETGATE(idt[IRQ_OFFSET + 3], 0, GD_KT, &ENTRY_IRQ3, 0);
7   SETGATE(idt[IRQ_OFFSET + 4], 0, GD_KT, &ENTRY_IRQ4, 0);
8   SETGATE(idt[IRQ_OFFSET + 5], 0, GD_KT, &ENTRY_IRQ5, 0);
9   SETGATE(idt[IRQ_OFFSET + 6], 0, GD_KT, &ENTRY_IRQ6, 0);
10  SETGATE(idt[IRQ_OFFSET + 7], 0, GD_KT, &ENTRY_IRQ7, 0);
11  SETGATE(idt[IRQ_OFFSET + 8], 0, GD_KT, &ENTRY_IRQ8, 0);
12  SETGATE(idt[IRQ_OFFSET + 9], 0, GD_KT, &ENTRY_IRQ9, 0);
13  SETGATE(idt[IRQ_OFFSET + 10], 0, GD_KT, &ENTRY_IRQ10, 0);
14  SETGATE(idt[IRQ_OFFSET + 11], 0, GD_KT, &ENTRY_IRQ11, 0);
15  SETGATE(idt[IRQ_OFFSET + 12], 0, GD_KT, &ENTRY_IRQ12, 0);
16  SETGATE(idt[IRQ_OFFSET + 13], 0, GD_KT, &ENTRY_IRQ13, 0);
17  SETGATE(idt[IRQ_OFFSET + 14], 0, GD_KT, &ENTRY_IRQ14, 0);
18  SETGATE(idt[IRQ_OFFSET + 15], 0, GD_KT, &ENTRY_IRQ15, 0);

```

In `env_alloc()`, we initialize the value of `%eflags` for user environments to be `FL_IF` so that user environments always run with interrupt enabled.

```

1   // Enable interrupts while in user mode.
2   // LAB 4: Your code here.
3   e->env_tf.tf_eflags = FL_IF;

```

By the way, since we design JOS so that interrupt is always disabled in kernel, all entries in IDT are interrupt gates which will clear IF bit when being passed through.

And now we set the IF bit when CPU is idle in `sched_halt()` to allow it to be awoken by interrupts.

```

1   // Reset stack pointer, enable interrupts and then halt.
2   asm volatile (
3       "movl $0, %%ebp\n"
4       "movl %0, %%esp\n"
5       "pushl $0\n"
6       "pushl $0\n"
7       // Uncomment the following line after completing exercise 13
8       "sti\n"
9       "1:\n"
10      "hlt\n"
11      "jmp 1b\n"
12      : : "a" (thiscpu->cpu_ts.ts_esp0));

```

Handling Clock Interrupts

Exercise 14

When interrupt with vector `IRQ_OFFSET + TIMER` is delivered, we transfer the control to `sched_yield()`. But remember to acknowledge the interrupt or it won't be generated again.

```

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

```

Inter-Process communication (IPC)

Implementing IPC

Exercise 15

In `sys_ipc_try_send()`, we try to query the receiver's `env_ipc_recving` to see if it's ready for an IPC

```
1 struct Env* e;
2 int r;
3
4 if ((r = envid2env(envid, &e, false)))
5 {
6     return r;
7 }
8
9 if(!e->env_ipc_recving)
10 {
11     return -E_IPC_NOT_RECV;
12 }
```

Then if both sender and receiver are willing to transfer a page, we check if the arguments are valid and then perform the same operations as in `sys_page_map()` (we don't call it directly since in `sys_page_map()` the caller must pass the permission check in `envid2env()` for both the source and destination while IPC can happen between any environments)

```
1 // only tries to transfer a page when they're both willing
2 if ((intptr_t)(srcva) < UTOP && (intptr_t)(e->env_ipc_dstva) < UTOP)
3 {
4     if ((intptr_t)(srcva) % PGSIZE)
5     {
6         return -E_INVALID;
7     }
8     if ((~perm & PTE_P) || (~perm & PTE_U) || (perm & ~PTE_SYSCALL))
9     {
10         return -E_INVALID;
11     }
12
13     // we don't call sys_page_map() to do this since it
14     // have more strict permission request
15     struct PageInfo* p;
16     pte_t *pte;
17     p = page_lookup(curenv->env_pgdir, srcva, &pte);
18     if (p == NULL)
19     {
20         return -E_INVALID;
21     }
22     if ((perm & PTE_W) && (~(*pte) & PTE_W))
23     {
24         return -E_INVALID;
25     }
26     if ((r = page_insert(e->env_pgdir, p, e->env_ipc_dstva, perm)))
27     {
```

```

28         return r;
29     }
30 }
31 else // now all checks have been passed at this position
32 {
33     e->env_ipc_perm = 0;
34 }

```

Then we fill in the IPC relevant field in the environment receiving and mark it as runnable again. The receiver's `%eax` stored in its `env_tf` is modified to 0 to act as the return value of `sys_ipc_recv()` called by the receiver.

```

1     e->env_ipc_recving = 0;
2     e->env_ipc_from = curenv->env_id;
3     e->env_ipc_value = value;
4     e->env_ipc_perm = perm;
5
6     e->env_status = ENV_RUNNABLE; // again we don't use sys_env_set_status for
    permission reason
7     e->env_tf.tf_regs.reg_eax = 0; // return 0 for sys_ipc_recv()
8
9     return 0;

```

In `sys_ipc_rev()`, we first check the validity of `dstva` and set `env_ipc_recving` true. Then we mark current environment as not runnable and give up CPU.

```

1     static int
2     sys_ipc_recv(void *dstva)
3     {
4         // LAB 4: Your code here.
5         if ((intptr_t)(dstva) < UTOP && (intptr_t)(dstva) % PGSIZE)
6         {
7             return -E_INVAL;
8         }
9         curenv->env_ipc_recving = true;
10        curenv->env_ipc_dstva = dstva;
11
12        curenv->env_status = ENV_NOT_RUNNABLE;
13        sched_yield();
14
15        // panic("sys_ipc_recv not implemented");
16
17        return 0; // the function actually doesn't return here
18    }

```

For the wrapper function, if the argument `pg` is set as `NULL` to indicate the unwilling of transferring a page in the IPC, we translate it as `UTOP` as an address `>= UTOP` is with the same semantic in the system calls. In `ipc_send()`, we keep trying to send the message and give up CPU after each failure because of `E_IPC_NOT_RECV` as a nice environment.

```

1     void
2     ipc_send(env_id_t to_env, uint32_t val, void *pg, int perm)
3     {

```

```

4 // LAB 4: Your code here.
5 int r;
6
7 // an address >= UTOP means no receiving page in sys_ipc_rcv()
8 pg = pg == NULL ? (void *)UTOP : pg;
9 do
10 {
11     r = sys_ipc_try_send(to_env, val, pg, perm);
12     if (r)
13     {
14         if (r != -E_IPC_NOT_RECV)
15         {
16             panic("ipc_send: %e", r);
17         }
18         else
19         {
20             sys_yield();
21         }
22     }
23 }
24 while(r);
25
26 // panic("ipc_send not implemented");
27 }
28
29 int32_t
30 ipc_rcv(envid_t *from_env_store, void *pg, int *perm_store)
31 {
32     // LAB 4: Your code here.
33     // an address >= UTOP means no receiving page in sys_ipc_rcv()
34     int r;
35     pg = pg == NULL ? (void *)UTOP : pg;
36     if ((r = sys_ipc_rcv(pg)))
37     {
38         if (from_env_store != NULL) *from_env_store = 0;
39         if (perm_store != NULL) *perm_store = 0;
40         return r;
41     }
42     if (from_env_store != NULL) *from_env_store = thisenv->env_ipc_from;
43     if (perm_store != NULL) *perm_store = thisenv->env_ipc_perm;
44     return thisenv->env_ipc_value;
45
46     // panic("ipc_rcv not implemented");
47     // return 0;
48 }

```

Challenge!

`ipc_send()` have to loop because `sys_ipc_try_send()` would fail when the receiver is not receiving or some other environment get to send to it before. This is an error recoverable in term of time. For example if the scheduler let a sender runs `sys_ipc_try_send()` before the receiver calls `sys_ipc_rcv()`, the first `sys_ipc_try_send()` will fail while the second one will success. That's we have to loop to wait for the receiver.

We can modify the IPC mechanism of the kernel to help environments to handle this. The idea is that we pause the sender if the receiver is not ready and record it in the receiver's waiting queue, and each time the receiver calls `sys_ipc_rcv()` we try to do IPC with the first environment in the queue (actually a last-in-first-out here because of my implementation) and wake the sender up on whether a success or failure, returning from `sys_ipc_try_send()` without the need of further loops.

We add two pointers in the environment structure since an environment can wait for a receiver to call `sys_ipc_rcv()` as some other environment are waiting for it to call `sys_ipc_rcv()`, which is the case in `user/primes`

```
1 struct Env *env_ipc_queue; // the head of IPC waiting queue (this is receiver)
2 struct Env *env_ipc_next; // next waiting environment in the same waiting queue (this is sender)
```

They are initialized in `env_alloc()`

```
1 // clean the IPC waiting queue pointers
2 e->env_ipc_queue = NULL;
3 e->env_ipc_next = NULL;
```

The changes in `syscall.c` are explained in the comments.

Note that it's ok to share the variables to store sending and receiving variables since one environment can not be a sender and a receiver simultaneous, although it can be the sender while some others are trying to send message to it as explained above.

```
1 // handle the IPC to dst from src (the head of dst's waiting queue)
2 // contains much of the original version of sys_ipc_try_send()
3 // can be called both from the sender or the receiver when
4 // (1) receiver calls sys_ipc_rcv() when some environments are waiting to send
5 // (2) sender calls sys_ipc_try_send() when the receiver is ready to receiver
6 // see sys_ipc_try_send() for possible errors and more information
7 static int
8 handle_ipc(struct Env* dst)
9 {
10     int r;
11
12     // pop the front of the waiting queue
13     struct Env* src = dst->env_ipc_queue;
14     assert(src != NULL);
15     dst->env_ipc_queue = src->env_ipc_next;
16
17     // restore the arguments from IPC relevant field
18     void *srcva = src->env_ipc_dstva;
19     uint32_t value = src->env_ipc_value;
20     int perm = src->env_ipc_perm;
21
22     // only tries to transfer a page when they're both willing
23     if ((intptr_t)(srcva) < UTOP && (intptr_t)(dst->env_ipc_dstva) < UTOP)
```

```

24     {
25         if ((intptr_t)(srcva) % PGSIZE)
26         {
27             r = -E_INVALID;
28             goto ret;
29         }
30         if ((~perm & PTE_P) || (~perm & PTE_U) || (perm & ~PTE_SYSCALL))
31         {
32             r = -E_INVALID;
33             goto ret;
34         }
35
36         // we don't call sys_page_map() to do this since it
37         // have more strict permission request
38         struct PageInfo* p;
39         pte_t *pte;
40         p = page_lookup(src->env_pgdir, srcva, &pte);
41         if (p == NULL)
42         {
43             r = -E_INVALID;
44             goto ret;
45         }
46         if ((perm & PTE_W) && (~(*pte) & PTE_W))
47         {
48             r = -E_INVALID;
49             goto ret;
50         }
51         if ((r = page_insert(dst->env_pgdir, p, dst->env_ipc_dstva, perm)))
52         {
53             goto ret;
54         }
55     }
56     else // now all checks have been passed at this position
57     {
58         dst->env_ipc_perm = 0;
59     }
60
61     dst->env_ipc_recving = 0;
62     dst->env_ipc_from = src->env_id;
63     dst->env_ipc_value = value;
64     dst->env_ipc_perm = perm;
65
66     r = 0;
67 ret:
68     // store return value in sender's or receiver's %eax
69     // in case they're sleeping
70     if (src->env_status == ENV_NOT_RUNNABLE)
71     {
72         src->env_status = ENV_RUNNABLE;
73         src->env_tf.tf_regs.reg_eax = r;
74     }
75     if (dst->env_status == ENV_NOT_RUNNABLE && !r) // receiver only wake up
on success
76     {
77         dst->env_status = ENV_RUNNABLE;
78         dst->env_tf.tf_regs.reg_eax = r;

```

```

79     }
80     // cprintf("handl_ipc(): from %x to %x, value = %d, retval = %d\n",
81     //         src->env_id, dst->env_id, value, r);
82     return r;
83 }
84
85 static int
86 sys_ipc_try_send(env_id_t env_id, uint32_t value, void *srcva, unsigned perm)
87 {
88     // LAB 4: Your code here.
89     struct Env* e;
90     int r;
91
92     if ((r = env_id2env(env_id, &e, false)))
93     {
94         return r;
95     }
96
97     // add current environment to the head of waiting queue of receiving
    environemnt
98     curenv->env_ipc_next = e->env_ipc_queue;
99     e->env_ipc_queue = curenv;
100    // store the arguments in IPC relevant field
101    curenv->env_ipc_value = value;
102    curenv->env_ipc_dstva = srcva;
103    curenv->env_ipc_perm = perm;
104
105    if(!e->env_ipc_recving)
106    {
107        // return -E_IPC_NOT_RECV;
108
109        // give up CPU if receiver isn't ready
110        // instead of return -E_IPC_NOT_RECV
111        curenv->env_status = ENV_NOT_RUNNABLE;
112        sched_yield();
113    }
114    // otherwise do the IPC
115    return handle_ipc(e);
116
117    // panic("sys_ipc_try_send not implemented");
118 }
119
120 static int
121 sys_ipc_recv(void *dstva)
122 {
123     // LAB 4: Your code here.
124     int r;
125
126     if ((intptr_t)(dstva) < UTOP && (intptr_t)(dstva) % PGSIZE)
127     {
128         return -E_INVAL;
129     }
130     curenv->env_ipc_recving = true;
131     curenv->env_ipc_dstva = dstva;
132
133     // travel IPC waiting queue (loop because some might fail)

```

```

134     while (curenv->env_ipc_queue != NULL)
135     {
136         r = handle_ipc(curenv);
137         if (!r) return r;
138     }
139
140     // no valid waiting environment, give up the CPU
141     curenv->env_status = ENV_NOT_RUNNABLE;
142     sched_yield();
143
144     // panic("sys_ipc_recv not implemented");
145
146     return 0; // the function actually doesn't return here
147 }

```

Now we can implement an `ipc_send()` without loops.

```

1  void
2  ipc_send(envid_t to_env, uint32_t val, void *pg, int perm)
3  {
4      // LAB 4: Your code here.
5      int r;
6
7      // an address >= UTOP means no receiving page in sys_ipc_recv()
8      pg = pg == NULL ? (void *)UTOP : pg;
9      r = sys_ipc_try_send(to_env, val, pg, perm);
10     if (r)
11     {
12         panic("ipc_send: %e", r);
13     }
14
15     // panic("ipc_send not implemented");
16 }

```

Reference

[1] IA32 volume 32, section 6.2.1