# MIT 6.828: Operating System Engineering

Lab Report
Lab 4: Preemptive Multitasking

CSNLP
csnlp16@126.com
https://csnlp.github.io/

# Contents

```
1    struct list {
2      int data;
3      struct list *next;
4    };
5
6    struct list *list = 0;
7
8    void
9    insert(int data)
10   {
11     struct list *l;
12
13     l = malloc(sizeof *l);
14     l->data = data;
15     l->next = list;
16     list = l;
17   }
```

Figure 1: An example of race condition

# 1 Introduction

## 1.1 Locking

A lock provides mutual exclusion, ensuring that only one CPU at a time can hold the lock. If a lock is associated with each shared data item. When using a given item, the code always holds the associated lock. In this way, we can be sure that the item is used from only one CPU at a time.

Give an example to show a *race condition*:

There will be a race between two CPU in line 16 since they both use the data item **list**.

To avoid race, the common way is to use a lock.

```
6     struct list *list = 0;
      struct lock listlock;
7
8     void
9     insert(int data)
10    {
11      struct list *l;
12      l = malloc(sizeof *l);
13      l->data = data;
14

        acquire(&listlock);
15      l->next = list;
16      list = l;
        release(&listlock);
17    }
```

Figure 2: Avoid race condition with a lock

## 1.2 The Lock Implementation in xv6

There are two types of locks in xv6: **spin-lock** and **sleep-lock**.

### 1.2.1 Spin-lock in xv6

The important field in spin-lock is **locked**, a word that is zero when the lock is available and non-zero when it is held. To acquire a lock, xv6 execute the following code:

```
21    void
22    acquire(struct spinlock *lk)
23    {
24      for(;;) {
25        if(!lk->locked) {
26          lk->locked = 1;
27          break;
28        }
29      }
30    }
```

Figure 3: Spin-lock: race in its code

This code cannot ensure mutual exclusion. When two CPU simultaneously reach line 25, they both grab the lock by executing line 26. In other words, they both acquire the lock. So, this implementation of **acquire** has its own race condition. To solve this problem, **line 25 and line 26 must execute in one *atomic* step.**

## 1.3  How To Use Locks

A hard part about using locks is deciding how many locks to use and which data and invariants each lock protects. There are a few basic principles:

1. Any time a variable can be **written** by one CPU at the same time that another CPU can **read or write** to it, a lock should be introduced to keep the two operations from overlapping.

2. Keep in mind that locks are used to protect **invariants**: if an invariant involves multiple memory locations, typically all of them need to be protected by a single lock to ensure the invariant is maintained.

## 1.4  Deadlock and Lock Ordering

### 1.4.1  How Deadlock Happens?

Let's say two code paths need lock A and lock B.

1. The first code path acquires locks in the order {A then B}.

2. The second code path acquires locks in the order {B then A}.

3. If thread T1 executes code path 1 while T2 executes code path 2 concurrently. T1 has acquired lock A while T2 has acquired lock B. They both face the situation that the other thread holds the needed lock and won't release it until its acquire returns.

The above situation will cause deadlock.

## 1.5  Sleep Locks

## 1.6  Limitation of Locks

## 1.7  New source files in LAB 4

1. **kern/cpu.h**: Kernel-private definitions for multiprocessor support.

2. **kern/mpconfig.c**: Code to read the multiprocessor configuration.

3. **kern/lapic**: Kernel code driving the local APIC unit in each processor.

4. **kern/mpentry.S**: Assembly-language entry code for non-boot CPUs.

5. **kern/spinlock.h**: Kernel-private definitions for spin locks, including the big kernel lock.

6. **kern/spinlock.c**: Kernel code implementing spin locks.

7. **kern/sched.c**: Code skeleton of the scheduler that you are about to implement.

# 2   Part A: Multiprocessor Support and Cooperative Multitasking

## 2.1   Multiprocessor Support

In this part, we expect that JOS can support *symmetric multiprocessing* (SMP), a multiprocessor model in which all CPUs have equivalent access to system resources such as memory and I/O buses.

During the boot process, these functionally identical CPUs can be classified into two types:

1. Bootstrap processor (BSP): responsible for initializing the system and for booting the OS.

2. Application processors (APs): be activated by the BSP only after the OS is up and running.

In an SMP system, each CPU has an accompanying local APIC (LAPIC) unit. The LAPIC units are responsible for:

1. Delivering interrupts throughout the system.

2. Providing its connected CPU with a unique identifier.

Which basic functionality of the LAPIC unit we will use:

1. APIC ID: to tell which CPU on which our code is current running, see in **cpunum()**.

2. Sending the STARTUP interprocessor interrupt(IPI) from the BSP to the APs to bring up other CPUs, see in **lapic_startap**

3. The LAPIC's built-in timer, see in **apic_init()**.

### 2.1.1   Exercise 1

*Implement **mmio_map_region** in **kern/pmap.c**. To see how this is used, look at the beginning of **lapic_init** in **kern/lapic.c**. You'll have to do the next exercise, too, before the tests for **mmio_map_region** will run.*

Let's firstly check **lapic_init()** in **kern/lapic.c**

```
void
lapic_init(void)
{
  if (!lapicaddr)
    return;

  // lapicaddr is the physical address of the LAPIC's 4K MMIO
  // region.  Map it in to virtual memory so we can access it.
  lapic = mmio_map_region(lapicaddr, 4096);

  // Enable local APIC; set spurious interrupt vector.
  lapicw(SVR, ENABLE | (IRQ_OFFSET + IRQ_SPURIOUS));

  // The timer repeatedly counts down at bus frequency
```

6

```
15      // from lapic[TICR] and then issues an interrupt.
16      // If we cared more about precise timekeeping,
17      // TICR would be calibrated using an external time source.
18      lapicw(TDCR, X1);
19      lapicw(TIMER, PERIODIC | (IRQ_OFFSET + IRQ_TIMER));
20      lapicw(TICR, 10000000);
21
22      // Leave LINT0 of the BSP enabled so that it can get
23      // interrupts from the 8259A chip.
24      //
25      // According to Intel MP Specification, the BIOS should initialize
26      // BSP's local APIC in Virtual Wire Mode, in which 8259A's
27      // INTR is virtually connected to BSP's LINTIN0. In this mode,
28      // we do not need to program the IOAPIC.
29      if (thiscpu != bootcpu)
30        lapicw(LINT0, MASKED);
31
32      // Disable NMI (LINT1) on all CPUs
33      lapicw(LINT1, MASKED);
34
35      // Disable performance counter overflow interrupts
36      // on machines that provide that interrupt entry.
37      if (((lapic[VER]>>16) & 0xFF) >= 4)
38        lapicw(PCINT, MASKED);
39
40      // Map error interrupt to IRQ_ERROR.
41      lapicw(ERROR, IRQ_OFFSET + IRQ_ERROR);
42
43      // Clear error status register (requires back-to-back writes).
44      lapicw(ESR, 0);
45      lapicw(ESR, 0);
46
47      // Ack any outstanding interrupts.
48      lapicw(EOI, 0);
49
50      // Send an Init Level De-Assert to synchronize arbitration ID's.
51      lapicw(ICRHI, 0);
52      lapicw(ICRLO, BCAST | INIT | LEVEL);
53      while(lapic[ICRLO] & DELIVS)
54        ;
55
56      // Enable interrupts on the APIC (but not on the processor).
57      lapicw(TPR, 0);
58  }
```

It calls **mmio_map_region()**

```
1           // lapicaddr is the physical address of the LAPIC's 4K MMIO
2           // region.  Map it in to virtual memory so we can access it.
3           lapic = mmio_map_region(lapicaddr, 4096);
```

Let's see the implementation of **mmio_map_region**:

```
1  void *
2  mmio_map_region(physaddr_t pa, size_t size)
3  {
4          static uintptr_t base = MMIOBASE;
5          void *ret = (void *)base;
6          size = ROUNDUP(size, PGSIZE);
```

```
 7          if(base + size > MMIOLIM || base + size < base)
 8                  panic("mmio_map_region fail: overflow");
 9          boot_map_region(kern_pgdir, base, size, pa, PTE_W | PTE_PCD | PTE_PWT);
10          base += size;
11          return ret;
12  }
```

### 2.1.2   Application Processor Bootstrap

1. **Information Collection**: Before booting up APs, the BSP should first collect information about the multiprocessor system, such as the total number of CPUs, their APIC IDs and the MMIO address of the LAPIC unit. The **mp_init()** function in **kern/mpconfig.c** retrieves this information by reading the MP configuration table that resides in the BIOS's region of memory. Let's see the **mp_init()** in **kern/mpconfig.c**

```
 1  void
 2  mp_init(void)
 3  {
 4    struct mp *mp;
 5    struct mpconf *conf;
 6    struct mpproc *proc;
 7    uint8_t *p;
 8    unsigned int i;
 9
10    bootcpu = &cpus[0];
11    if ((conf = mpconfig(&mp)) == 0)
12      return;
13    ismp = 1;
14    lapicaddr = conf->lapicaddr;
15
16    for (p = conf->entries, i = 0; i < conf->entry; i++) {
17      switch (*p) {
18      case MPPROC:
19        proc = (struct mpproc *)p;
20        if (proc->flags & MPPROC_BOOT)
21          bootcpu = &cpus[ncpu];
22        if (ncpu < NCPU) {
23          cpus[ncpu].cpu_id = ncpu;
24          ncpu++;
25        } else {
26          cprintf("SMP: too many CPUs, CPU %d disabled\n",
27            proc->apicid);
28        }
29        p += sizeof(struct mpproc);
30        continue;
31      case MPBUS:
32      case MPIOAPIC:
33      case MPIOINTR:
34      case MPLINTR:
35        p += 8;
36        continue;
37      default:
38        cprintf("mpinit: unknown config type %x\n", *p);
39        ismp = 0;
40        i = conf->entry;
```

```
41      }
42    }
43
44    bootcpu->cpu_status = CPU_STARTED;
45    if (!ismp) {
46      // Didn't like what we found; fall back to no MP.
47      ncpu = 1;
48      lapicaddr = 0;
49      cprintf("SMP: configuration not found, SMP disabled\n");
50      return;
51    }
52    cprintf("SMP: CPU %d found %d CPU(s)\n", bootcpu->cpu_id,  ncpu);
53
54    if (mp->imcrp) {
55      // [MP 3.2.6.1] If the hardware implements PIC mode,
56      // switch to getting interrupts from the LAPIC.
57      cprintf("SMP: Setting IMCR to switch from PIC mode to symmetric I/O mode\n")
      ↪ ;
58      outb(0x22, 0x70);   // Select IMCR
59      outb(0x23, inb(0x23) | 1);  // Mask external interrupts.
60    }
61 }
```

2. **AP bootstrap process**: The **boot_aps()** function (in **kern/init.c**) drives the AP bootstrap process. APs start in real mode, much like how the bootloader started in **boot/boot.S**, so **boot_aps()** copies the AP entry code (**kern/mpentry.S**) to a memory location that is addressable in the real mode. Unlike with the bootloader, we have some control over where the AP will start executing code; we copy the entry code to **0x7000 (MPENTRY_PADDR)**, but any unused, page-aligned physical address below 640KB would work.

3. **AP's activation**: **boot_aps()** activates APs one after another, by sending **STARTUP IPIs** to the LAPIC unit of the corresponding AP, along with an initial **CS:IP** address at which the AP should start running its entry code (**MPENTRY_PADDR** in our case). The entry code in **kern/mpentry.S** is quite similar to that of **boot/-boot.S**. After some brief setup, it puts the AP into protected mode with paging enabled, and then calls the C setup routine **mp_main()** (also in **kern/init.c**). **boot_aps()** waits for the AP to signal a **CPU_STARTED** flag in **cpu_status** field of its **struct CpuInfo** before going on to wake up the next one.

```
1  // Start the non-boot (AP) processors.
2  static void
3  boot_aps(void)
4  {
5          extern unsigned char mpentry_start[], mpentry_end[];
6          void *code;
7          struct CpuInfo *c;
8
9          // Write entry code to unused memory at MPENTRY_PADDR
10         code = KADDR(MPENTRY_PADDR);
11         memmove(code, mpentry_start, mpentry_end - mpentry_start);
12
13         // Boot each AP one at a time
14         for (c = cpus; c < cpus + ncpu; c++) {
15                 if (c == cpus + cpunum())  // We've started already.
```

9

```
16                    continue;
17
18              // Tell mpentry.S what stack to use
19              mpentry_kstack = percpu_kstacks[c - cpus] + KSTKSIZE;
20              // Start the CPU at mpentry_start
21              lapic_startap(c->cpu_id, PADDR(code));
22              // Wait for the CPU to finish some basic setup in mp_main()
23              while(c->cpu_status != CPU_STARTED)
24                       ;
25         }
26 }
```

The **boot_aps()** is about:

- **mpentry_start** and **mpentry_end** is the start and end of the virtual address of the assembly code: **kern/mpentry.S** in memory.

- **memmove** is to move the assembly code **kern/mpentry.S** to physical address: **MPENTRY_PADDR**.

- **mpentry_kstack** is stack of the boot loader of the selected AP.

- **lapic_startap**: start the CPU at the **MPENTRY_PADDR**, notice that at this time **MPENTRY_PADDR** is already the start address of **kern/mpentry.S**.

### 2.1.3 Exercise 2

**Exercise 2.** Read `boot_aps()` and `mp_main()` in `kern/init.c`, and the assembly code in `kern/mpentry.S`. Make sure you understand the control flow transfer during the bootstrap of APs. Then modify your implementation of `page_init()` in `kern/pmap.c` to avoid adding the page at MPENTRY_PADDR to the free list, so that we can safely copy and run AP bootstrap code at that physical address. Your code should pass the updated `check_page_free_list()` test (but might fail the updated `check_kern_pgdir()` test, which we will fix soon).

**Question**

1. Compare `kern/mpentry.S` side by side with `boot/boot.S`. Bearing in mind that `kern/mpentry.S` is compiled and linked to run above KERNBASE just like everything else in the kernel, what is the purpose of macro MPBOOTPHYS? Why is it necessary in `kern/mpentry.S` but not in `boot/boot.S`? In other words, what could go wrong if it were omitted in `kern/mpentry.S`?
Hint: recall the differences between the link address and the load address that we have discussed in Lab 1.

Figure 4: Exercise 2

```
1  // Setup code for APs
2  void
3  mp_main(void)
4  {
5          // We are in high EIP now, safe to switch to kern_pgdir
6          lcr3(PADDR(kern_pgdir));
7          cprintf("SMP: CPU %d starting\n", cpunum());
8
9          lapic_init();
10         env_init_percpu();
11         trap_init_percpu();
12         xchg(&thiscpu->cpu_status, CPU_STARTED); // tell boot_aps() we're up
13
14         // Now that we have finished some basic setup, call sched_yield()
15         // to start running processes on this CPU.  But make sure that
16         // only one CPU can enter the scheduler at a time!
17         //
18         // Your code here:
19
20         // Remove this after you finish Exercise 6
21         for (;;);
22  }
```

Listing 1: map_main in kern/init.c

### 2.1.3.1   The Control Flow

The control flow transfer during the bootstrap of APs can be summarized as follows:

1. **mem_init()**: Lab 2 memory management initialization.

2. **env_init()** and **trap_init()**: Lab 3 user environment initialization.

3. **mp_init()** and **lapic_init()**: Lab 4 multiprocessor initialization.

4. **pic_init()**: multitasking initialization.

5. **lock_kernel()**: acquire the kernel before waking up APs.

6. **boot_aps()**: starting non-boot CPUs.

### 2.1.4   Per-CPU State and Initialization

When writing a multiprocessor OS, it's important to distinguish between:

- per-CPU state that is private to each processor.

- Global state that the whole system shares.

Here is the per-CPU state that we should be aware of:

- **per-CPU kernel stack:** Since multiple CPUs can trap into the kernel simultaneously, we need a separate kernel stack for each processor to prevent them from interfering with each other's execution. See the actual structure in **inc/memlayout.h**

11

```c
#ifndef JOS_INC_MEMLAYOUT_H
#define JOS_INC_MEMLAYOUT_H

#ifndef __ASSEMBLER__
#include <inc/types.h>
#include <inc/mmu.h>
#endif /* not __ASSEMBLER__ */

/*
 * This file contains definitions for memory management in our OS,
 * which are relevant to both the kernel and user-mode software.
 */

// Global descriptor numbers
#define GD_KT     0x08     // kernel text
#define GD_KD     0x10     // kernel data
#define GD_UT     0x18     // user text
#define GD_UD     0x20     // user data
#define GD_TSS0   0x28     // Task segment selector for CPU 0

/*
 * Virtual memory map:                               Permissions
 *                                                   kernel/user
 *
 *    4 Gig -------->  +------------------------------+
 *                     |                              | RW/--
 *                     ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~
 *                     :              .               :
 *                     :              .               :
 *                     :              .               :
 *                     |~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~| RW/--
 *                     |                              | RW/--
 *                     |    Remapped Physical Memory  | RW/--
 *                     |                              | RW/--
 *    KERNBASE, ---->  +------------------------------+ 0xf0000000      --+
 *    KSTACKTOP        |     CPU0's Kernel Stack      | RW/--  KSTKSIZE   |
 *                     | - - - - - - - - - - - - - - -|                   |
 *                     |      Invalid Memory (*)      | --/--  KSTKGAP    |
 *                     +------------------------------+                   |
 *                     |     CPU1's Kernel Stack      | RW/--  KSTKSIZE   |
 *                     | - - - - - - - - - - - - - - -|                 PTSIZE
 *                     |      Invalid Memory (*)      | --/--  KSTKGAP    |
 *                     +------------------------------+                   |
 *                     :              .               :                   |
 *                     :              .               :                   |
 *    MMIOLIM ------>  +------------------------------+ 0xefc00000      --+
 *                     |       Memory-mapped I/O      | RW/--  PTSIZE
 * ULIM, MMIOBASE -->  +------------------------------+ 0xef800000
 *                     |  Cur. Page Table (User R-)   | R-/R-  PTSIZE
 *    UVPT      ---->  +------------------------------+ 0xef400000
 *                     |          RO PAGES            | R-/R-  PTSIZE
 *    UPAGES    ---->  +------------------------------+ 0xef000000
 *                     |          RO ENVS             | R-/R-  PTSIZE
 * UTOP,UENVS ------>  +------------------------------+ 0xeec00000
 * UXSTACKTOP -/       |     User Exception Stack     | RW/RW  PGSIZE
 *                     +------------------------------+ 0xeebff000
 *                     |       Empty Memory (*)       | --/--  PGSIZE
 *    USTACKTOP  --->  +------------------------------+ 0xeebfe000
```

```
 59  *                          |      Normal User Stack      | RW/RW  PGSIZE
 60  *                          +------------------------------+ 0xeebfd000
 61  *                          |                              |
 62  *                          |                              |
 63  *                          ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~
 64  *                          .                              .
 65  *                          .                              .
 66  *                          .                              .
 67  *                          |~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~|
 68  *                          |       Program Data & Heap     |
 69  *    UTEXT -------->  +------------------------------+ 0x00800000
 70  *    PFTEMP ------->  |       Empty Memory (*)       |        PTSIZE
 71  *                          |                              |
 72  *    UTEMP -------->  +------------------------------+ 0x00400000      --+
 73  *                          |       Empty Memory (*)       |                |
 74  *                          | - - - - - - - - - - - - - - -|                |
 75  *                          |  User STAB Data (optional)   |             PTSIZE
 76  *    USTABDATA ---->  +------------------------------+ 0x00200000      |
 77  *                          |       Empty Memory (*)       |                |
 78  *    0 ------------>  +------------------------------+                  --+
 79  *
 80  * (*) Note: The kernel ensures that "Invalid Memory" is *never* mapped.
 81  *     "Empty Memory" is normally unmapped, but user programs may map pages
 82  *     there if desired.  JOS user programs map pages temporarily at UTEMP.
 83  */
 84
 85
 86 // All physical memory mapped at this address
 87 #define KERNBASE   0xF0000000
 88
 89 // At IOPHYSMEM (640K) there is a 384K hole for I/O.  From the kernel,
 90 // IOPHYSMEM can be addressed at KERNBASE + IOPHYSMEM.  The hole ends
 91 // at physical address EXTPHYSMEM.
 92 #define IOPHYSMEM 0x0A0000
 93 #define EXTPHYSMEM  0x100000
 94
 95 // Kernel stack.
 96 #define KSTACKTOP KERNBASE
 97 #define KSTKSIZE  (8*PGSIZE)      // size of a kernel stack
 98 #define KSTKGAP   (8*PGSIZE)      // size of a kernel stack guard
 99
100 // Memory-mapped IO.
101 #define MMIOLIM   (KSTACKTOP - PTSIZE)
102 #define MMIOBASE  (MMIOLIM - PTSIZE)
103
104 #define ULIM     (MMIOBASE)
105
106 /*
107  * User read-only mappings! Anything below here til UTOP are readonly to user.
108  * They are global pages mapped in at env allocation time.
109  */
110
111 // User read-only virtual page table (see 'uvpt' below)
112 #define UVPT    (ULIM - PTSIZE)
113 // Read-only copies of the Page structures
114 #define UPAGES    (UVPT - PTSIZE)
115 // Read-only copies of the global env structures
116 #define UENVS    (UPAGES - PTSIZE)
```

```c
/*
 * Top of user VM. User can manipulate VA from UTOP-1 and down!
 */

// Top of user-accessible VM
#define UTOP      UENVS
// Top of one-page user exception stack
#define UXSTACKTOP  UTOP
// Next page left invalid to guard against exception stack overflow; then:
// Top of normal user stack
#define USTACKTOP (UTOP - 2*PGSIZE)

// Where user programs generally begin
#define UTEXT    (2*PTSIZE)

// Used for temporary page mappings.  Typed 'void*' for convenience
#define UTEMP    ((void*) PTSIZE)
// Used for temporary page mappings for the user page-fault handler
// (should not conflict with other temporary page mappings)
#define PFTEMP    (UTEMP + PTSIZE - PGSIZE)
// The location of the user-level STABS data structure
#define USTABDATA (PTSIZE / 2)

// Physical address of startup code for non-boot CPUs (APs)
#define MPENTRY_PADDR 0x7000

#ifndef __ASSEMBLER__

typedef uint32_t pte_t;
typedef uint32_t pde_t;

#if JOS_USER
/*
 * The page directory entry corresponding to the virtual address range
 * [UVPT, UVPT + PTSIZE) points to the page directory itself.  Thus, the page
 * directory is treated as a page table as well as a page directory.
 *
 * One result of treating the page directory as a page table is that all PTEs
 * can be accessed through a "virtual page table" at virtual address UVPT (to
 * which uvpt is set in lib/entry.S).  The PTE for page number N is stored in
 * uvpt[N].  (It's worth drawing a diagram of this!)
 *
 * A second consequence is that the contents of the current page directory
 * will always be available at virtual address (UVPT + (UVPT >> PGSHIFT)), to
 * which uvpd is set in lib/entry.S.
 */
extern volatile pte_t uvpt[];     // VA of "virtual page table"
extern volatile pde_t uvpd[];     // VA of current page directory
#endif

/*
 * Page descriptor structures, mapped at UPAGES.
 * Read/write to the kernel, read-only to user programs.
 *
 * Each struct PageInfo stores metadata for one physical page.
 * Is it NOT the physical page itself, but there is a one-to-one
 * correspondence between physical pages and struct PageInfo's.
```

```
175  * You can map a struct PageInfo * to the corresponding physical address
176  * with page2pa() in kern/pmap.h.
177  */
178 struct PageInfo {
179   // Next page on the free list.
180   struct PageInfo *pp_link;
181
182   // pp_ref is the count of pointers (usually in page table entries)
183   // to this page, for pages allocated using page_alloc.
184   // Pages allocated at boot time using pmap.c's
185   // boot_alloc do not have valid reference count fields.
186
187   uint16_t pp_ref;
188 };
189
190 #endif /* !__ASSEMBLER__ */
191 #endif /* !JOS_INC_MEMLAYOUT_H */
```

Listing 2: inc/memlayout.h

- **Per-CPU TSS and TSS descriptor:** A per-CPU task state segment (TSS) is needed to specify each CPU's kernel lives. The TSS for CPU $i$ is stored in **cpu[i].cpu_ts**, and the corresponding TSS descriptor is defined in the GDT entry **gdt[(GD_TSS0 ¿¿ 3) + i]**.

- **Per-CPU current environment pointer:** Because each CPU can run different user process simultaneously, we redefine they symbol **curenv** to refer to **cpus[cpunum()].cpu_env** (or **thiscpu¿cpu_env**), which points to the environment currently executing on the current CPU. **Per-CPU system registers:** All registers are private to CPU. Therefore, instructions initializing these registers (such as **lcr3(), ltr(), lgdt(), lidt()**) must be executed once on each CPU. Functions **env_init_percpu()** and **trap_init_percpu()** are defined for this purpose.

### 2.1.5 Exercise 3

**Exercise 3.** Modify `mem_init_mp()` (in `kern/pmap.c`) to map per–CPU stacks starting at `KSTACKTOP`, as shown in `inc/memlayout.h`. The size of each stack is `KSTKSIZE` bytes plus `KSTKGAP` bytes of unmapped guard pages. Your code should pass the new check in `check_kern_pgdir()`.

Figure 5: Exercise 3

Let's firstly see hints in **mem_init_mp()**:

```
1 // Modify mappings in kern_pgdir to support SMP
2 //    - Map the per-CPU stacks in the region [KSTACKTOP-PTSIZE, KSTACKTOP)
3 //
4 static void
5 mem_init_mp(void)
```

```
6  {
7          // Map per-CPU stacks starting at KSTACKTOP, for up to 'NCPU' CPUs.
8          //
9          // For CPU i, use the physical memory that 'percpu_kstacks[i]' refers
10         // to as its kernel stack. CPU i's kernel stack grows down from virtual
11         // address kstacktop_i = KSTACKTOP - i * (KSTKSIZE + KSTKGAP), and is
12         // divided into two pieces, just like the single stack you set up in
13         // mem_init:
14         //      * [kstacktop_i - KSTKSIZE, kstacktop_i)
15         //            -- backed by physical memory
16         //      * [kstacktop_i - (KSTKSIZE + KSTKGAP), kstacktop_i - KSTKSIZE)
17         //            -- not backed; so if the kernel overflows its stack,
18         //               it will fault rather than overwrite another CPU's stack.
19         //               Known as a "guard page".
20         //      Permissions: kernel RW, user NONE
21         //
22         // LAB 4: Your code here:
23
24 }
```

See its implementation:

```
1  static void
2  mem_init_mp(void)
3  {
4          uintptr_t kstacktop_i;
5          int i;
6          for(i=0; i < NCPU; i++) {
7                  kstacktop_i = KSTACKTOP - i * (KSTKSIZE + KSTKGAP);
8                  boot_map_region(kern_pgdir,
9                                  kstacktop_i - KSTKSIZE,
10                                 ROUNDUP(KSTKSIZE, PGSIZE),
11                                 PADDR(&percpu_kstacks[i]),
12                                 PTE_W | PTE_P);
13         }
14
15 }
```

### 2.1.6  Exercise 4

**Exercise 4.** The code in `trap_init_percpu()` (`kern/trap.c`) initializes the TSS and TSS descriptor for the BSP. It worked in Lab 3, but is incorrect when running on other CPUs. Change the code so that it can work on all CPUs. (Note: your new code should not use the global `ts` variable any more.)

Figure 6: Exercise 4

Firstly, see its hints:

```
1  // Initialize and load the per-CPU TSS and IDT
2  void
3  trap_init_percpu(void)
```

```
4  {
5    // The example code here sets up the Task State Segment (TSS) and
6    // the TSS descriptor for CPU 0. But it is incorrect if we are
7    // running on other CPUs because each CPU has its own kernel stack.
8    // Fix the code so that it works for all CPUs.
9    //
10   // Hints:
11   //   - The macro "thiscpu" always refers to the current CPU's
12   //     struct CpuInfo;
13   //   - The ID of the current CPU is given by cpunum() or
14   //     thiscpu->cpu_id;
15   //   - Use "thiscpu->cpu_ts" as the TSS for the current CPU,
16   //     rather than the global "ts" variable;
17   //   - Use gdt[(GD_TSS0 >> 3) + i] for CPU i's TSS descriptor;
18   //   - You mapped the per-CPU kernel stacks in mem_init_mp()
19   //   - Initialize cpu_ts.ts_iomb to prevent unauthorized environments
20   //     from doing IO (0 is not the correct value!)
21   //
22   // ltr sets a 'busy' flag in the TSS selector, so if you
23   // accidentally load the same TSS on more than one CPU, you'll
24   // get a triple fault.  If you set up an individual CPU's TSS
25   // wrong, you may not get a fault until you try to return from
26   // user space on that CPU.
27   //
28   // LAB 4: Your code here:
29
30   // Setup a TSS so that we get the right stack
31   // when we trap to the kernel.
32   ts.ts_esp0 = KSTACKTOP;
33   ts.ts_ss0 = GD_KD;
34   ts.ts_iomb = sizeof(struct Taskstate);
35
36   // Initialize the TSS slot of the gdt.
37   gdt[GD_TSS0 >> 3] = SEG16(STS_T32A, (uint32_t) (&ts),
38           sizeof(struct Taskstate) - 1, 0);
39   gdt[GD_TSS0 >> 3].sd_s = 0;
40
41   // Load the TSS selector (like other segment selectors, the
42   // bottom three bits are special; we leave them 0)
43   ltr(GD_TSS0);
44
45   // Load the IDT
46   lidt(&idt_pd);
47  }
```

See the SMP version:

```
1  void
2  trap_init_percpu(void)
3  {
4          thiscpu->cpu_ts.ts_esp0 = KSTACKTOP - cpunum() * (KSTKGAP + KSTKSIZE);
5          thiscpu->cpu_ts.ts_ss0 = GD_KD;
6
7          // Initialize the TSS slot of the gdt.
8          gdt[(GD_TSS0 >> 3) + cpunum()] = SEG16(STS_T32A,
9                                        (uint32_t)(&thiscpu->cpu_ts),
   ↪                              sizeof(struct Taskstate)-1,
10                                        0);
11         gdt[(GD_TSS0 >> 3) + cpunum()].sd_s = 0;
```

```
12          // Load the TSS selector (like other segment selectors, the
13          // bottom three bits are special; we leave them 0)
14          ltr(GD_TSS0 + sizeof(struct Segdesc) * cpunum());
15
16          // Load the IDT
17          lidt(&idt_pd);
18  }
```

**It JUST RUN AGAIN AND AGAIN!!!!!**

### 2.1.7  Locking

Before letting the AP get any further, we need to firstly address race conditions when multiple CPUs run kernel code simultaneously. The simplest way to achieve this is to use a *big kernel lock*. The big kernel lock is a single lock that is held whenever an environment enters **kernel mode**, and released when the environment returns to **user mode**. In this model:

1. Environments in user mode can run concurrently on any available CPUs.

2. No more that one environment can run in kernel mode. Any other environments trying to enter kernel mode are forced to wait.

This kernel lock is declared in **kern/spinlock.h**, namely **kernel_lock**. It also provides **lock_kernel()** and **unlock_kernel()** to acquire and release the lock.

This kernel lock can be applied in the following locations:

1. **i386_init()**: acquire the lock before the BSP wakes up the other CPUs.

2. **mp_main()**: acquire the lock after initializing the AP. Then, call **sched_yield()** to start running environments on this AP.

3. **trap()**: acquire the lock when trapped from the user mode. To determine whether the trap is from the user mode or kernel mode, check the low bits of the **tf_cs**.

4. **env_run()**: release the lock right before switching to user mode.

### 2.1.8  Exercise 5

*Apply the big kernel lock as described above, by calling **lock_kernel()** and **unlock_kernel()** at the proper locations.*

1. **i386_init()** in **kern/init.c**

```
1
2   void
3   i386_init(void)
4   {
5     // Initialize the console.
6     // Can't call cprintf until after we do this!
7     cons_init();
8
9     cprintf("6828 decimal is %o octal!\n", 6828);
10
11    // Lab 2 memory management initialization functions
12    mem_init();
```

18

```
13
14    // Lab 3 user environment initialization functions
15    env_init();
16    trap_init();
17
18    // Lab 4 multiprocessor initialization functions
19    mp_init();
20    lapic_init();
21
22    // Lab 4 multitasking initialization functions
23    pic_init();
24
25    // Acquire the big kernel lock before waking up APs
26    // Your code here:
27    lock_kernel();
28
29    // Starting non-boot CPUs
30    boot_aps();
31
32 #if defined(TEST)
33    // Don't touch -- used by grading script!
34    ENV_CREATE(TEST, ENV_TYPE_USER);
35 #else
36    // Touch all you want.
37    ENV_CREATE(user_primes, ENV_TYPE_USER);
38 #endif // TEST*
39
40    // Schedule and run the first user environment!
41    sched_yield();
42 }
```

2. **mp_main()** in kern/init.c

```
1 // Setup code for APs
2 void
3 mp_main(void)
4 {
5         // We are in high EIP now, safe to switch to kern_pgdir
6         lcr3(PADDR(kern_pgdir));
7         cprintf("SMP: CPU %d starting\n", cpunum());
8
9         lapic_init();
10        env_init_percpu();
11        trap_init_percpu();
12        xchg(&thiscpu->cpu_status, CPU_STARTED); // tell boot_aps() we're up
13
14        // Now that we have finished some basic setup, call sched_yield()
15        // to start running processes on this CPU.  But make sure that
16        // only one CPU can enter the scheduler at a time!
17        //
18        // Your code here:
19        lock_kernel();
20        sched_yield();
21        // Remove this after you finish Exercise 6
22        for (;;);
23 }
```

3. **trap()** in kern/trap.c

19

```
1  void
2  trap(struct Trapframe *tf)
3  {
4    // The environment may have set DF and some versions
5    // of GCC rely on DF being clear
6    asm volatile("cld" ::: "cc");
7
8    // Halt the CPU if some other CPU has called panic()
9    extern char *panicstr;
10   if (panicstr)
11     asm volatile("hlt");
12
13   // Re-acqurie the big kernel lock if we were halted in
14   // sched_yield()
15   if (xchg(&thiscpu->cpu_status, CPU_STARTED) == CPU_HALTED)
16     lock_kernel();
17   // Check that interrupts are disabled.  If this assertion
18   // fails, DO NOT be tempted to fix it by inserting a "cli" in
19   // the interrupt path.
20   assert(!(read_eflags() & FL_IF));
21
22   if ((tf->tf_cs & 3) == 3) {
23     // Trapped from user mode.
24     // Acquire the big kernel lock before doing any
25     // serious kernel work.
26     // LAB 4: Your code here.
27     lock_kernel();
28     assert(curenv);
29
30     // Garbage collect if current enviroment is a zombie
31     if (curenv->env_status == ENV_DYING) {
32       env_free(curenv);
33       curenv = NULL;
34       sched_yield();
35     }
36
37     // Copy trap frame (which is currently on the stack)
38     // into 'curenv->env_tf', so that running the environment
39     // will restart at the trap point.
40     curenv->env_tf = *tf;
41     // The trapframe on the stack should be ignored from here on.
42     tf = &curenv->env_tf;
43   }
44
45   // Record that tf is the last real trapframe so
46   // print_trapframe can print some additional information.
47   last_tf = tf;
48
49   // Dispatch based on what type of trap occurred
50   trap_dispatch(tf);
51
52   // If we made it to this point, then no other environment was
53   // scheduled, so we should return to the current environment
54   // if doing so makes sense.
55   if (curenv && curenv->env_status == ENV_RUNNING)
56     env_run(curenv);
57   else
58     sched_yield();
```

```
59 }
```

4. **env_run()** in kern/env.c

```
1  // Context switch from curenv to env e.
2  // Note: if this is the first call to env_run, curenv is NULL.
3  //
4  // This function does not return.
5  //
6  void
7  env_run(struct Env *e)
8  {
9    // Step 1: If this is a context switch (a new environment is running):
10   //      1. Set the current environment (if any) back to
11   //         ENV_RUNNABLE if it is ENV_RUNNING (think about
12   //         what other states it can be in),
13   //      2. Set 'curenv' to the new environment,
14   //      3. Set its status to ENV_RUNNING,
15   //      4. Update its 'env_runs' counter,
16   //      5. Use lcr3() to switch to its address space.
17   // Step 2: Use env_pop_tf() to restore the environment's
18   //      registers and drop into user mode in the
19   //      environment.
20
21   // Hint: This function loads the new environment's state from
22   //  e->env_tf.  Go back through the code you wrote above
23   //  and make sure you have set the relevant parts of
24   //  e->env_tf to sensible values.
25
26   // LAB 3: Your code here.
27   // Step 1: change the curenv's env_status from ING to ABLE.
28   if(curenv != NULL && curenv->env_status == ENV_RUNNING)
29     curenv->env_status = ENV_RUNNABLE;
30
31   // Step 2: switch to the new env
32   curenv = e;
33   // Step 3: set status
34   curenv->env_status = ENV_RUNNING;
35   // Step 4: update env_runs counter
36   curenv->env_runs++;
37   // Step 5: use lcr3() to switch to its address space
38   lcr3(PADDR(curenv->env_pgdir));
39
40   // LAB 4: unlock_kernel()
41   unlock_kernel();
42   // Step 6: use env_pop_tf() to restore the env's register
43   //      and drop into user mode in the environment.
44   env_pop_tf(&curenv->env_tf);
45
46
47   panic("env_run not yet implemented");
48 }
```

## 2.2 Round-Robin Scheduling

Change the JOS kernel so that it can alternate between multiple environments in *round-robin* [1] fashion:

- **sched_yield()** in **kern/sched.c**: responsible for selecting a new environment to run. It searches sequentially through the **envs[]** array in circular fashion, picks the first environment it finds with a status of **ENV_RUNNABLE** (see in inc/env.h) and calls **env_run()** to jump to the selected environment. Notice that **sched_yield()** must **NEVER** run the same environment on two CPUs at the same time.

- Implement a new system call **sys_yield()**: by which user environment can call to invoke kernel's **sched_yield()** function and thereby **voluntarily give up** the CPU to a different environment.

### 2.2.1 Exercise 6

**Exercise 6.** Implement round–robin scheduling in `sched_yield()` as described above. Don't forget to modify `syscall()` to dispatch `sys_yield()`.

Make sure to invoke `sched_yield()` in `mp_main`.

Modify `kern/init.c` to create three (or more!) environments that all run the program `user/yield.c`.

Run `make qemu`. You should see the environments switch back and forth between each other five times before terminating, like below.

Test also with several CPUS: `make qemu CPUS=2`.

```
...
Hello, I am environment 00001000.
Hello, I am environment 00001001.
Hello, I am environment 00001002.
Back in environment 00001000, iteration 0.
Back in environment 00001001, iteration 0.
Back in environment 00001002, iteration 0.
Back in environment 00001000, iteration 1.
Back in environment 00001001, iteration 1.
Back in environment 00001002, iteration 1.
...
```

After the `yield` programs exit, there will be no runnable environment in the system, the scheduler should invoke the JOS kernel monitor. If any of this does not happen, then fix your code before proceeding.

Figure 7: Exercise 6

---

[1]Round-robin is one of the algorithms employed by process and network schedulers in computing. As the term is generally used, time slices are assigned to each process in equal portions and in circular order, handling all processes without priority.

### 2.2.1.1 sched_yield()

Firstly, see its hints:

```
1  // Choose a user environment to run and run it.
2  void
3  sched_yield(void)
4  {
5          struct Env *idle;
6
7          // Implement simple round-robin scheduling.
8          //
9          // Search through 'envs' for an ENV_RUNNABLE environment in
10         // circular fashion starting just after the env this CPU was
11         // last running.  Switch to the first such environment found.
12         //
13         // If no envs are runnable, but the environment previously
14         // running on this CPU is still ENV_RUNNING, it's okay to
15         // choose that environment.
16         //
17         // Never choose an environment that's currently running on
18         // another CPU (env_status == ENV_RUNNING). If there are
19         // no runnable environments, simply drop through to the code
20         // below to halt the cpu.
21
22         // LAB 4: Your code here.
23
24         // sched_halt never returns
25         sched_halt();
26 }
```

The implementation:

```
1  // Choose a user environment to run and run it.
2  void
3  sched_yield(void)
4  {
5    struct Env *idle;
6
7    // LAB 4: Your code here.
8    uint32_t i, j, start;
9    struct Env *foundenv;
10   foundenv = NULL;
11   // get the current env running on this CPU.
12   struct Env *cpu_curenv = thiscpu->cpu_env;
13   start = (cpu_curenv != NULL) ? ENVX(cpu_curenv->env_id) : 0;
14   // get the env with the least priority.
15   for(i=0; i < NENV; i++) {
16     j = (start + i) % NENV;
17       if (envs[j].env_status == ENV_RUNNABLE) {
18         if (foundenv == NULL)
19           foundenv = &envs[j];
20           break;
21       }
22   }
23   // check if the kernel has found available CPU if not use the curenv
24   if ((cpu_curenv && cpu_curenv->env_status == ENV_RUNNING) && (foundenv == NULL)) {
25     env_run(cpu_curenv);
26     return;
```

```
27      }
28      if(foundenv) {
29          env_run(foundenv);
30          return;
31      }
32
33
34      // sched_halt never returns
35      sched_halt();
36  }
```

### 2.2.1.2 Modify syscall() to Dispatch sys_yield()

Add a case **SYS_yield()**.

```
1  // Dispatches to the correct kernel function, passing the arguments.
2  int32_t
3  syscall(uint32_t syscallno, uint32_t a1, uint32_t a2, uint32_t a3, uint32_t a4,
   ↪ uint32_t a5)
4  {
5          // Call the function corresponding to the 'syscallno' parameter.
6          // Return any appropriate return value.
7          // LAB 3: Your code here.
8
9
10         //panic("syscall not implemented");
11
12         switch (syscallno) {
13         case SYS_cputs:
14                 sys_cputs((char *)a1, a2);
15                 return 0;
16         case SYS_cgetc:
17                 return sys_cgetc();
18         case SYS_getenvid:
19                 return sys_getenvid();
20         case SYS_env_destroy:
21                 return sys_env_destroy(a1);
22         case SYS_yield:
23                 sys_yield();
24                 return 0;
25         default:
26                 return -E_INVAL;
27         }
28 }
```

### 2.2.1.3 Modify i386_init() in kern/init.c to Add More Environments

```
1  #if defined(TEST)
2          // Don't touch -- used by grading script!
3          ENV_CREATE(TEST, ENV_TYPE_USER);
4  #else
5          // Touch all you want.
6          //ENV_CREATE(user_primes, ENV_TYPE_USER);
7          ENV_CREATE(user_yield, ENV_TYPE_USER);
8          ENV_CREATE(user_yield, ENV_TYPE_USER);
9          ENV_CREATE(user_yield, ENV_TYPE_USER);
```

```
10  #endif // TEST*
11
12          // Schedule and run the first user environment!
13          sched_yield();
```

The result is as follows:

```
1   + cc kern/init.c
2   + ld obj/kern/kernel
3   + mk obj/kern/kernel.img
4   sed "s/localhost:1234/localhost:26000/" < .gdbinit.tmpl > .gdbinit
5   qemu-system-i386 -drive file=obj/kern/kernel.img,index=0,media=disk,format=raw -serial
        ↪   mon:stdio -gdb tcp::26000 -D qemu.log -smp 1
6   6828 decimal is 15254 octal!
7   Physical memory: 131072K available, base = 640K, extended = 130432K
8   check_page_free_list() succeeded!
9   check_page_alloc() succeeded!
10  check_page() succeeded!
11  check_kern_pgdir() succeeded!
12  check_page_free_list() succeeded!
13  check_page_installed_pgdir() succeeded!
14  SMP: CPU 0 found 1 CPU(s)
15  enabled interrupts: 1 2
16  [00000000] new env 00001000
17  [00000000] new env 00001001
18  [00000000] new env 00001002
19  Hello, I am environment 00001000.
20  Hello, I am environment 00001001.
21  Hello, I am environment 00001002.
22  Back in environment 00001000, iteration 0.
23  Back in environment 00001001, iteration 0.
24  Back in environment 00001002, iteration 0.
25  Back in environment 00001000, iteration 1.
26  Back in environment 00001001, iteration 1.
27  Back in environment 00001002, iteration 1.
28  Back in environment 00001000, iteration 2.
29  Back in environment 00001001, iteration 2.
30  Back in environment 00001002, iteration 2.
31  Back in environment 00001000, iteration 3.
32  Back in environment 00001001, iteration 3.
33  Back in environment 00001002, iteration 3.
34  Back in environment 00001000, iteration 4.
35  All done in environment 00001000.
36  [00001000] exiting gracefully
37  [00001000] free env 00001000
38  Back in environment 00001001, iteration 4.
39  All done in environment 00001001.
40  [00001001] exiting gracefully
41  [00001001] free env 00001001
42  Back in environment 00001002, iteration 4.
43  All done in environment 00001002.
44  [00001002] exiting gracefully
45  [00001002] free env 00001002
46  No runnable environments in the system!
47  Welcome to the JOS kernel monitor!
48  Type 'help' for a list of commands.
```

## 2.3 System Calls for Environment Creation

Now implement the necessary JOS system calls to allow user environments to create and start other new user environments. Unix provides the **fork()** system call as its process creation primitive.

- **fork()** copies the entire address space of the calling process (the parent) to create a new process (the child).

- In the parent, the **fork()** returns the child's process IDs.

- In the child, the **fork()** returns 0.

- Both the parent and child have its private address space.

Implement a more primitive system call to create a new environment.

- **sys_exofork()**: Create a new environment.

- **sys_env_set_status()**: Set status.

- **sys_page_alloc**: Allocate page.

- **sys_page_map**: Copy a page mapping.

- **sys_page_unmap**: Unmap a page.

### 2.3.1 Exercise 7

**Exercise 7.** Implement the system calls described above in `kern/syscall.c` and make sure `syscall()` calls them. You will need to use various functions in `kern/pmap.c` and `kern/env.c`, particularly `envid2env()`. For now, whenever you call `envid2env()`, pass 1 in the `checkperm` parameter. Be sure you check for any invalid system call arguments, returning `-E_INVAL` in that case. Test your JOS kernel with `user/dumbfork` and make sure it works before proceeding.

Figure 8: Exercise 7

### 2.3.2 The Implementation of sys_exofork()

Let's see the hints:

```
// Allocate a new environment.
// Returns envid of new environment, or < 0 on error.  Errors are:
//      -E_NO_FREE_ENV if no free environment is available.
//      -E_NO_MEM on memory exhaustion.
static envid_t
sys_exofork(void)
{
        // Create the new environment with env_alloc(), from kern/env.c.
```

26

```
 9          // It should be left as env_alloc created it, except that
10          // status is set to ENV_NOT_RUNNABLE, and the register set is copied
11          // from the current environment -- but tweaked so sys_exofork
12          // will appear to return 0.
13
14          // LAB 4: Your code here.
15          panic("sys_exofork not implemented");
16 }
```

See the implementation:

```
 1 static envid_t
 2 sys_exofork(void)
 3 {
 4          // Create the new environment with env_alloc(), from kern/env.c.
 5          // It should be left as env_alloc created it, except that
 6          // status is set to ENV_NOT_RUNNABLE, and the register set is copied
 7          // from the current environment -- but tweaked so sys_exofork
 8          // will appear to return 0.
 9
10          // LAB 4: Your code here.
11          struct Env *newenv;
12          int32_t ret = 0;
13          if ((ret = env_alloc(&newenv, curenv->env_id)) < 0) {
14                  return ret;
15          }
16          newenv->env_status = ENV_NOT_RUNNABLE;
17          newenv->env_tf = curenv->env_tf;
18          newenv->env_tf.tf_regs.reg_eax = 0;
19          return newenv->env_id;
20          //panic("sys_exofork not implemented");
21 }
```

### 2.3.3 The Implementation of sys_env_set_status()

Sets the status of a specified environment to **ENV_RUNNABLE** or **ENV_NOT_RUNNABLE**.
This system call is typically used to mark a new environment ready to run, once its address space and register state has been fully initialized.

### 2.3.3.1 See The Initial Hints:

```
 1 // Set envid's env_status to status, which must be ENV_RUNNABLE
 2 // or ENV_NOT_RUNNABLE.
 3 //
 4 // Returns 0 on success, < 0 on error.  Errors are:
 5 //      -E_BAD_ENV if environment envid doesn't currently exist,
 6 //              or the caller doesn't have permission to change envid.
 7 //      -E_INVAL if status is not a valid status for an environment.
 8 static int
 9 sys_env_set_status(envid_t envid, int status)
10 {
11          // Hint: Use the 'envid2env' function from kern/env.c to translate an
12          // envid to a struct Env.
13          // You should set envid2env's third argument to 1, which will
14          // check whether the current environment has permission to set
15          // envid's status.
```

```
16
17          // LAB 4: Your code here.
18          panic("sys_env_set_status not implemented");
19 }
```

### 2.3.3.2    See the implementation

```
1  static int
2  sys_env_set_status(envid_t envid, int status)
3  {
4          // Hint: Use the 'envid2env' function from kern/env.c to translate an
5          // envid to a struct Env.
6          // You should set envid2env's third argument to 1, which will
7          // check whether the current environment has permission to set
8          // envid's status.
9
10          // LAB 4: Your code here.
11          struct Env *env;
12          // Check whether the current environment is a valid status
13          if(status != ENV_RUNNABLE && status != ENV_NOT_RUNNABLE)
14                  return -E_INVAL;
15          // Check whether the current environment has permission to set
16          // envid's status
17          if (envid2env(envid, &env, 1) < 0)
18                  return -E_BAD_ENV;
19          env->env_status = status;
20          return 0;
21          //panic("sys_env_set_status not implemented");
22 }
```

### 2.3.3.3    The envid2env() Function in kern/env.c

```
1  //
2  // Converts an envid to an env pointer.
3  // If checkperm is set, the specified environment must be either the
4  // current environment or an immediate child of the current environment.
5  //
6  // RETURNS
7  //   0 on success, -E_BAD_ENV on error.
8  //   On success, sets *env_store to the environment.
9  //   On error, sets *env_store to NULL.
10 //
11 int
12 envid2env(envid_t envid, struct Env **env_store, bool checkperm)
13 {
14   struct Env *e;
15
16   // If envid is zero, return the current environment.
17   if (envid == 0) {
18     *env_store = curenv;
19     return 0;
20   }
21
22   // Look up the Env structure via the index part of the envid,
23   // then check the env_id field in that struct Env
24   // to ensure that the envid is not stale
25   // (i.e., does not refer to a _previous_ environment
26   // that used the same slot in the envs[] array).
```

```
27    e = &envs[ENVX(envid)];
28    if (e->env_status == ENV_FREE || e->env_id != envid) {
29      *env_store = 0;
30      return -E_BAD_ENV;
31    }
32
33    // Check that the calling environment has legitimate permission
34    // to manipulate the specified environment.
35    // If checkperm is set, the specified environment
36    // must be either the current environment
37    // or an immediate child of the current environment.
38    if (checkperm && e != curenv && e->env_parent_id != curenv->env_id) {
39      *env_store = 0;
40      return -E_BAD_ENV;
41    }
42
43    *env_store = e;
44    return 0;
45 }
```

## 2.3.4   The Implementation of sys_page_alloc()

Allocate a page of physical memory and maps it at a given virtual address in a given environment's address space.

### 2.3.4.1   The hints

```
1  // Allocate a page of memory and map it at 'va' with permission
2  // 'perm' in the address space of 'envid'.
3  // The page's contents are set to 0.
4  // If a page is already mapped at 'va', that page is unmapped as a
5  // side effect.
6  //
7  // perm -- PTE_U | PTE_P must be set, PTE_AVAIL | PTE_W may or may not be set,
8  //         but no other bits may be set.  See PTE_SYSCALL in inc/mmu.h.
9  //
10 // Return 0 on success, < 0 on error.  Errors are:
11 //      -E_BAD_ENV if environment envid doesn't currently exist,
12 //              or the caller doesn't have permission to change envid.
13 //      -E_INVAL if va >= UTOP, or va is not page-aligned.
14 //      -E_INVAL if perm is inappropriate (see above).
15 //      -E_NO_MEM if there's no memory to allocate the new page,
16 //              or to allocate any necessary page tables.
17 static int
18 sys_page_alloc(envid_t envid, void *va, int perm)
19 {
20         // Hint: This function is a wrapper around page_alloc() and
21         //   page_insert() from kern/pmap.c.
22         //   Most of the new code you write should be to check the
23         //   parameters for correctness.
24         //   If page_insert() fails, remember to free the page you
25         //   allocated!
26
27         // LAB 4: Your code here.
28         panic("sys_page_alloc not implemented");
29 }
```

### 2.3.4.2 See the implementation

```
static int
sys_page_alloc(envid_t envid, void *va, int perm)
{
        // LAB 4: Your code here.
        struct Env *env;
        // check whether environment envid exists or not.
        if (envid2env(envid, &env, 1) < 0)
                return -E_BAD_ENV;
        // Check if va >= UTOP or va is not page-aligned.
        if((uintptr_t) va >= UTOP || PGOFF(va))
                return -E_INVAL;
        // Check if perm is inappropriate
        if((perm & PTE_SYSCALL) == 0)
                return -E_INVAL;
        if((perm & ~PTE_SYSCALL))
                return -E_INVAL;

        struct PageInfo *pp = page_alloc(ALLOC_ZERO);
        if(!pp)
                return -E_NO_MEM;
        if (page_insert(env->env_pgdir, pp, va, perm) < 0){
                // page_insert fails: remember free the page
                page_free(pp);
                return -E_NO_MEM;
        }

        // Success, return 0
        return 0;
}
```

## 2.3.5  The Implementation of sys_page_map()

**sys_page_map()** is used to copy a page mapping(**NOT the contents of a page**) from one environment's address space to another, leaving a memory sharing arrangement in place so that the new and the old mappings both refer to the same page of physical memory.

### 2.3.5.1  The Hints

### 2.3.5.2  Some hints from page_lookup()

```
//
// Return the page mapped at virtual address 'va'.
// If pte_store is not zero, then we store in it the address
// of the pte for this page.  This is used by page_remove and
// can be used to verify page permissions for syscall arguments,
// but should not be used by most callers.
//
// Return NULL if there is no page mapped at va.
//
// Hint: the TA solution uses pgdir_walk and pa2page.
//
struct PageInfo *
```

```
13  page_lookup(pde_t *pgdir, void *va, pte_t **pte_store)
14  {
15  }
```

### 2.3.5.3 Some hints from page_insert()

```
1   // Map the physical page 'pp' at virtual address 'va'.
2   // The permissions (the low 12 bits) of the page table entry
3   // should be set to 'perm|PTE_P'.
4   //
5   // Requirements
6   //    - If there is already a page mapped at 'va', it should be page_remove()d.
7   //    - If necessary, on demand, a page table should be allocated and inserted
8   //      into 'pgdir'.
9   //    - pp->pp_ref should be incremented if the insertion succeeds.
10  //    - The TLB must be invalidated if a page was formerly present at 'va'.
11  //
12  // Corner-case hint: Make sure to consider what happens when the same
13  // pp is re-inserted at the same virtual address in the same pgdir.
14  // However, try not to distinguish this case in your code, as this
15  // frequently leads to subtle bugs; there's an elegant way to handle
16  // everything in one code path.
17  //
18  // RETURNS:
19  //    0 on success
20  //    -E_NO_MEM, if page table couldn't be allocated
21  //
22  // Hint: The TA solution is implemented using pgdir_walk, page_remove,
23  // and page2pa.
24  //
25  int
26  page_insert(pde_t *pgdir, struct PageInfo *pp, void *va, int perm)
27  {
28  }
```

### 2.3.5.4 The Implementation of sys_page_map()

```
1   static int
2   sys_page_map(envid_t srcenvid, void *srcva,
3                envid_t dstenvid, void *dstva, int perm)
4   {
5       // LAB 4: Your code here.
6       struct Env *srcenv, *dstenv;
7       struct PageInfo *pp;
8       pte_t *pte;
9       // srcenvid or/and dstenvid doesn't exist.
10      if ((envid2env(srcenvid, &srcenv, 1) < 0) || (envid2env(dstenvid, &dstenv, 1)
    ↪ < 0)) {
11              return -E_BAD_ENV;
12      }
13      // Check srcva>= UTOP or srcva is not page-aligned
14      if ((uintptr_t)srcva >= UTOP || PGOFF(srcva))
15              return -E_INVAL;
16      // Check dstva>= UTOP or dstva is not page-aligned
17      if ((uintptr_t)dstva >= UTOP || PGOFF(dstva))
18              return -E_INVAL;
19      // Check the perm is inappropriate
20      if (((perm & PTE_SYSCALL) == 0)|| (perm & ~PTE_SYSCALL))
```

```
21            return -E_INVAL;
22        // Look up the page
23        pp = page_lookup(srcenv->env_pgdir, srcva, &pte);
24        if (pp == NULL)
25                return -E_INVAL;
26
27        // Check if (perm & PTE_W), but srcva is read-only in
28        // srcenvid's address space.
29        if ((perm & PTE_W) && ((*pte & PTE_W) == 0))
30                return -E_INVAL;
31        // Check if there's no memory in dstenv to allocate
32        // any necessary page tables.
33        if(page_insert(dstenv->env_pgdir, pp, dstva, perm) < 0)
34                return -E_NO_MEM;
35        return 0;
36        //panic("sys_page_map not implemented");
37 }
```

### 2.3.5.5    Summary

Its overall steps can be summarized as follow:

1. Check the parameters.

2. Use **page_lookup()** to find the physical page **pp** mapped by **srcva** in **srcenv**.

3. Map the **pp** to **dstva** in **dstenv**.

### 2.3.6    The Implementation of sys_page_unmap()

Unmap a page mapped at a given virtual address in a given environment.
    Let's firstly check the hints:

### 2.3.6.1    The Hints of sys_page_unmap()

```
1  // Unmap the page of memory at 'va' in the address space of 'envid'.
2  // If no page is mapped, the function silently succeeds.
3  //
4  // Return 0 on success, < 0 on error.  Errors are:
5  //      -E_BAD_ENV if environment envid doesn't currently exist,
6  //              or the caller doesn't have permission to change envid.
7  //      -E_INVAL if va >= UTOP, or va is not page-aligned.
8  static int
9  sys_page_unmap(envid_t envid, void *va)
10 {
11        // Hint: This function is a wrapper around page_remove().
12
13        // LAB 4: Your code here.
14        panic("sys_page_unmap not implemented");
15 }
```

The implementation is quite easy.

### 2.3.6.2    The Implementation

```
1 static int
2 sys_page_unmap(envid_t envid, void *va)
3 {
4       // Hint: This function is a wrapper around page_remove().
5
6       // LAB 4: Your code here.
7       struct Env *env;
8       if (envid2env(envid, &env, 1) < 0)
9               return -E_BAD_ENV;
10      if((uintptr_T)va >= UTOP || PGOFF(va))
11              return -E_INVAL;
12      page_remove(env->env_pgdir, va);
13      return 0;
14      //panic("sys_page_unmap not implemented");
15 }
```

### 2.3.7   Current syscall() in syscall.c

```
1 int32_t
2 syscall(uint32_t syscallno, uint32_t a1, uint32_t a2, uint32_t a3, uint32_t a4,
   ↪ uint32_t a5)
3 {
4       // Call the function corresponding to the 'syscallno' parameter.
5       // Return any appropriate return value.
6       // LAB 3: Your code here.
7
8
9       //panic("syscall not implemented");
10
11      switch (syscallno) {
12      case SYS_cputs:
13              sys_cputs((char *)a1, a2);
14              return 0;
15      case SYS_cgetc:
16              return sys_cgetc();
17      case SYS_getenvid:
18              return sys_getenvid();
19      case SYS_env_destroy:
20              return sys_env_destroy(a1);
21      case SYS_yield:
22              sys_yield();
23              return 0;
24      case SYS_page_alloc:
25              return sys_page_alloc((envid_t)a1, (void *)a2, (int)a3);
26      case SYS_page_map:
27              return sys_page_map((envid_t)a1, (void *)a2, (envid_t)a3, (void *)a4,
   ↪ (int)a5);
28      case SYS_page_unmap:
29              return sys_page_unmap((envid_t)a1, (void *)a2);
30      case SYS_exofork:
31              return sys_exofork();
32      case SYS_env_set_status:
33              return sys_env_set_status((envid_t)a1, (int)a2);
34      default:
35              return -E_INVAL;
36      }
37 }
```

# 3 Part B: Copy-on-Write Fork

## 3.1 The Motivation of Copy-on-Write Technique

The copying of the parent's space into the child is the most expensive part of the **fork()** operation. However, a call to **fork()** is frequently followed by a call to **exec()** in the child process, which replace the child's memory with a new program. Under this condition, the time spent copying the parent's address space is largely wasted because the child process will use very little of its memory before **exec()**

**Copy-on-Write Technique:** For this reason, later verion Unix took advantage of virtual memory hardware to allow the parent and child to **share** the memory mapped into the respective address space until one of the processes actually modifies it. This technique is know as *copy-on-write* technique. To do this,

1. On **fork()**, the kernel would copy the address space mappings from the parent to the child instead of the content of the mapped pages. At the same time, the kernel should mark the now-shared pages read-only.

2. When one of the two processes tries to write to one of these shared pages, the process take a page fault. At this point, the kernel realizes that the page was really a *virtual* or *copy-on-write* copy, and so the kernel makes a **new, private writable** copy of the page for the faulting process.

In this way, the contents of individual page aren't actually copied until they are actually written to. This optimization makes a **fork()** followed by an **exec()** in the child much cheaper.

## 3.2 User-level Page Fault Handling

### 3.2.1 Setting the Page Fault Handler

In order to handle its own page fault, a user environment will need to register a **page fault handler entrypoint** with the JOS kernel. This registration is done via **sys_env_set_pgfault_upcall()** system call.

#### 3.2.1.1 Exercise 8

**Exercise 8.** Implement the `sys_env_set_pgfault_upcall` system call. Be sure to enable permission checking when looking up the environment ID of the target environment, since this is a "dangerous" system call.

Figure 9: Exercise 8

#### 3.2.1.2 The hints in sys_env_set_pgfault_upcall()

```
1  // Set the page fault upcall for 'envid' by modifying the corresponding struct
2  // Env's 'env_pgfault_upcall' field.  When 'envid' causes a page fault, the
3  // kernel will push a fault record onto the exception stack, then branch to
4  // 'func'.
5  //
6  // Returns 0 on success, < 0 on error.  Errors are:
7  //      -E_BAD_ENV if environment envid doesn't currently exist,
8  //              or the caller doesn't have permission to change envid.
9  static int
10 sys_env_set_pgfault_upcall(envid_t envid, void *func)
11 {
12         // LAB 4: Your code here.
13         panic("sys_env_set_pgfault_upcall not implemented");
14 }
```

### 3.2.1.3   See The Implementation of sys_env_set_pgfault_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          // the environment envid doesn't currently exist
7          // or the caller has no permission to change envid
8          if(envid2env(envid, &env, 1) < 0)
9                  return -E_BAD_ENV;
10         // set the env_pgfault_upcall
11         env->env_pgfault_upcall = func;
12         return 0;
13
14         //panic("sys_env_set_pgfault_upcall not implemented");
15 }
```

## 3.2.2   Normal and Exception Stacks in User Environments

### 3.2.2.1   Normal User Stack and User Exception Stack

1. **Normal User Stack:**   During normal execution, a user environment in JOS will run on the normal user stack: its **ESP** register starts out pointing at **USTACK-TOP**, and the stack data it pushes resides on the page between **USTACKTOP-PGSIZE** and **USTACKTOP-1** inclusive.

2. **User Exception Stack:**   When a page fault occurs in user mode, however, the kernel will restart the user environment running a designated user-level page fault handler on a different stack, namely the **user exception stack**.

### 3.2.2.2   About JOS User Exception Stack

JOS user exception stack is also one page in size.

The JOS user exception stack is also one page in size, and its top is defined to be at virtual address UXSTACKTOP, so the valid bytes of the user exception stack are from UXSTACKTOP-PGSIZE through UXSTACKTOP-1 inclusive. While running on this exception stack, the user-level page fault handler can use JOS's regular system calls to map new pages or adjust mappings so as to fix whatever problem originally caused

the page fault. Then the user-level page fault handler returns, via an assembly language stub, to the faulting code on the original stack.

Each user environment that wants to support user-level page fault handling will need to allocate memory for its own exception stack, using the **sys_page_alloc()** system call introduced in part A.

### 3.2.3   Invoking the User Page Fault Handler

We will now need to change the page fault handling coder in **kern/trap.c** to handle page faults from user mode as follows.

1. If there is no page fault handler registered, the JOS kernel destroys the user environment with a message as before.

2. Otherwise

   - the kernel sets up a trap frame on the exception stack that looks like a **struct UTrapframe** from **inc/trap.h**:

```
                        <-- UXSTACKTOP
trap-time esp
trap-time eflags
trap-time eip
trap-time eax           start of struct PushRegs
trap-time ecx
trap-time edx
trap-time ebx
trap-time esp
trap-time ebp
trap-time esi
trap-time edi           end of struct PushRegs
tf_err (error code)
fault_va                <-- %esp when handler is run
```

Figure 10: Struct UTrapframe

   - The kernel then arranges for the user environment to resume execution with the page fault handler running on the exception stack with this stack frame; you must figure out how to make this happen. The fault_va is the virtual address that caused the page fault.

   - If the user environment is already running on the user exception stack when an exception occurs, then the page fault handler itself has faulted. In this case, you should start the new stack frame just under the current tf-¿tf_esp rather than at UXSTACKTOP. You should first push an empty 32-bit word, then a struct UTrapframe.

| Stack Type | Address |
|---|---|
| User Normal Stack | [UTEXT, USTACKTOP] |
| User Exception Stack | [UXSTACKTOP - PGSIZE, UXSTACKTOP] |
| Kernel Stack | [KSTACK-KSTKSIZE, KSTACKTOP] |

- To test whether tf-¿tf_esp is already on the user exception stack, check whether it is in the range between UXSTACKTOP-PGSIZE and UXSTACKTOP-1, inclusive.

NOW, the involved stack are:

### 3.2.4 Exercise 9

**Exercise 9.** Implement the code in `page_fault_handler` in `kern/trap.c` required to dispatch page faults to the user–mode handler. Be sure to take appropriate precautions when writing into the exception stack. (What happens if the user environment runs out of space on the exception stack?)

Figure 11: Exercise 9

Firstly, see its hints:

```
void
page_fault_handler(struct Trapframe *tf)
{
        uint32_t fault_va;

        // Read processor's CR2 register to find the faulting address
        fault_va = rcr2();

        // Handle kernel-mode page faults.

        // LAB 3: Your code here.
        if(tf->tf_cs && 0x01 == 0) {
                panic("page fault in kernel mode, fault address %d\n", fault_va);
        }

        // We've already handled kernel-mode exceptions, so if we get here,
        // the page fault happened in user mode.

        // Call the environment's page fault upcall, if one exists.  Set up a
        // page fault stack frame on the user exception stack (below
        // UXSTACKTOP), then branch to curenv->env_pgfault_upcall.
        //
        // The page fault upcall might cause another page fault, in which case
        // we branch to the page fault upcall recursively, pushing another
        // page fault stack frame on top of the user exception stack.
        //
        // It is convenient for our code which returns from a page fault
```

```
28          // (lib/pfentry.S) to have one word of scratch space at the top of the
29          // trap-time stack; it allows us to more easily restore the eip/esp. In
30          // the non-recursive case, we don't have to worry about this because
31          // the top of the regular user stack is free.  In the recursive case,
32          // this means we have to leave an extra word between the current top of
33          // the exception stack and the new stack frame because the exception
34          // stack _is_ the trap-time stack.
35          //
36          // If there's no page fault upcall, the environment didn't allocate a
37          // page for its exception stack or can't write to it, or the exception
38          // stack overflows, then destroy the environment that caused the fault.
39          // Note that the grade script assumes you will first check for the page
40          // fault upcall and print the "user fault va" message below if there is
41          // none.  The remaining three checks can be combined into a single test.
42          //
43          // Hints:
44          //   user_mem_assert() and env_run() are useful here.
45          //   To change what the user environment runs, modify 'curenv->env_tf'
46          //   (the 'tf' variable points at 'curenv->env_tf').

48          // LAB 4: Your code here.

50          // Destroy the environment that caused the fault.
51          cprintf("[%08x] user fault va %08x ip %08x\n",
52                  curenv->env_id, fault_va, tf->tf_eip);
53          print_trapframe(tf);
54          env_destroy(curenv);
55 }
```

### 3.2.4.1   Trapframe and UTrapframe

```
1  struct Trapframe {
2          struct PushRegs tf_regs;
3          uint16_t tf_es;
4          uint16_t tf_padding1;
5          uint16_t tf_ds;
6          uint16_t tf_padding2;
7          uint32_t tf_trapno;
8          /* below here defined by x86 hardware */
9          uint32_t tf_err;
10         uintptr_t tf_eip;
11         uint16_t tf_cs;
12         uint16_t tf_padding3;
13         uint32_t tf_eflags;
14         /* below here only when crossing rings, such as from user to kernel */
15         uintptr_t tf_esp;
16         uint16_t tf_ss;
17         uint16_t tf_padding4;
18 } __attribute__((packed));
19
20 struct UTrapframe {
21         /* information about the fault */
22         uint32_t utf_fault_va;  /* va for T_PGFLT, 0 otherwise */
23         uint32_t utf_err;
24         /* trap-time return state */
25         struct PushRegs utf_regs;
```

```
26        uintptr_t utf_eip;
27        uint32_t utf_eflags;
28        /* the trap-time stack to return to */
29        uintptr_t utf_esp;
30 } __attribute__((packed));
```

### 3.2.4.2  The Implementation of Exercise 9

```
 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 && 0x01 == 0) {
13     panic("page fault in kernel mode, fault address %d\n", fault_va);
14   }
15
16
17   // LAB 4: Your code here.
18
19   struct UTrapframe *utf;
20   if(curenv->env_pgfault_upcall) {
21     // tf->tf_esp in [UXSTACKTOP - PGSIZE, UXSTACKTOP-1]
22     if(UXSTACKTOP - PGSIZE <= tf->tf_esp && tf->tf_esp <= UXSTACKTOP - 1)
23       utf = (struct UTrapframe *)(tf->tf_esp - sizeof(struct UTrapframe) - 4);
24     else
25       utf = (struct UTrapframe *)(UXSTACKTOP - sizeof(struct UTrapframe));
26     // mem_assert
27     user_mem_assert(curenv, (void *)utf, sizeof(struct UTrapframe), PTE_U | PTE_W);
28     // set utf
29     utf->utf_fault_va = fault_va;
30     utf->utf_err = tf->tf_trapno;
31     utf->utf_eip = tf->tf_eip;
32     utf->utf_eflags = tf->tf_eflags;
33     utf->utf_esp = tf->tf_esp;
34     utf->utf_regs = tf->tf_regs;
35
36     // set tf
37     tf->tf_eip = (uint32_t)curenv->env_pgfault_upcall;
38     tf->tf_esp = (uint32_t)utf;
39     // run current env
40     env_run(curenv);
41   }else {
42
43   // Destroy the environment that caused the fault.
44   cprintf("[%08x] user fault va %08x ip %08x\n",
45     curenv->env_id, fault_va, tf->tf_eip);
46   print_trapframe(tf);
47   env_destroy(curenv);
48   }
49 }
```

### 3.2.5   User-mode Page Fault Entrypoint: Exercise 10 and Exercise 11

Next, you need to implement the assembly routine that will take care of calling the C page fault handler and resume execution at the original faulting instruction. This assembly routine is the handler that will be registered with the kernel using `sys_env_set_pgfault_upcall()`.

> **Exercise 10.** Implement the `_pgfault_upcall` routine in `lib/pfentry.s`. The interesting part is returning to the original point in the user code that caused the page fault. You'll return directly there, without going back through the kernel. The hard part is simultaneously switching stacks and re-loading the EIP.

Finally, you need to implement the C user library side of the user–level page fault handling mechanism.

> **Exercise 11.** Finish `set_pgfault_handler()` in `lib/pgfault.c`.

Figure 12: Exercise 10 and Exercise 11

### 3.2.5.1   See the hints in lib/pfentry.S

```
1  #include <inc/mmu.h>
2  #include <inc/memlayout.h>
3
4  // Page fault upcall entrypoint.
5
6  // This is where we ask the kernel to redirect us to whenever we cause
7  // a page fault in user space (see the call to sys_set_pgfault_handler
8  // in pgfault.c).
9  //
10 // When a page fault actually occurs, the kernel switches our ESP to
11 // point to the user exception stack if we're not already on the user
12 // exception stack, and then it pushes a UTrapframe onto our user
13 // exception stack:
14 //
15 //      trap-time esp
16 //      trap-time eflags
17 //      trap-time eip
18 //      utf_regs.reg_eax
19 //      ...
20 //      utf_regs.reg_esi
21 //      utf_regs.reg_edi
22 //      utf_err (error code)
23 //      utf_fault_va            <-- %esp
24 //
25 // If this is a recursive fault, the kernel will reserve for us a
26 // blank word above the trap-time esp for scratch work when we unwind
27 // the recursive call.
28 //
```

```
29  // We then have call up to the appropriate page fault handler in C
30  // code, pointed to by the global variable '_pgfault_handler'.
31  .text
32  .globl _pgfault_upcall
33  _pgfault_upcall:
34          // Call the C page fault handler.
35          pushl %esp                          // function argument: pointer to UTF
36          movl _pgfault_handler, %eax
37          call *%eax
38          addl $4, %esp                       // pop function argument
39
40          // Now the C page fault handler has returned and you must return
41          // to the trap time state.
42          // Push trap-time %eip onto the trap-time stack.
43          //
44          // Explanation:
45          //    We must prepare the trap-time stack for our eventual return to
46          //    re-execute the instruction that faulted.
47          //    Unfortunately, we can't return directly from the exception stack:
48          //    We can't call 'jmp', since that requires that we load the address
49          //    into a register, and all registers must have their trap-time
50          //    values after the return.
51          //    We can't call 'ret' from the exception stack either, since if we
52          //    did, %esp would have the wrong value.
53          //    So instead, we push the trap-time %eip onto the *trap-time* stack!
54          //    Below we'll switch to that stack and call 'ret', which will
55          //    restore %eip to its pre-fault value.
56          //
57          //    In the case of a recursive fault on the exception stack,
58          //    note that the word we're pushing now will fit in the
59          //    blank word that the kernel reserved for us.
60          //
61          // Throughout the remaining code, think carefully about what
62          // registers are available for intermediate calculations.  You
63          // may find that you have to rearrange your code in non-obvious
64          // ways as registers become unavailable as scratch space.
65        // LAB 4: Your code here.
66
67          // Restore the trap-time registers.  After you do this, you
68          // can no longer modify any general-purpose registers.
69          // LAB 4: Your code here.
70
71          // Restore eflags from the stack.  After you do this, you can
72          // no longer use arithmetic operations or anything else that
73          // modifies eflags.
74          // LAB 4: Your code here.
75
76          // Switch back to the adjusted trap-time stack.
77          // LAB 4: Your code here.
78
79          // Return to re-execute the instruction that faulted.
80          // LAB 4: Your code here.
```

#### 3.2.5.2 Implementation of the _pgfault_upcall routine in lib/pfentry.S

```
1  .text
2  .globl _pgfault_upcall
3  _pgfault_upcall:
```

```
4          // Call the C page fault handler.
5          pushl %esp                          // function argument: pointer to UTF
6          movl _pgfault_handler, %eax
7          call *%eax
8          addl $4, %esp                       // pop function argument
9          // LAB 4: Your code here.
10
11         // 48(%esp) = 48(%esp) - 4
12         movl 48(%esp), %ebp
13         subl $4, %ebp
14         movl %ebp, 48(%esp)
15
16         // %ebp = 40(%esp)
17         movl 40(%esp), %eax
18         movl %eax, (%ebp)
19
20         // Restore the trap-time registers.  After you do this, you
21         // can no longer modify any general-purpose registers.
22         // LAB 4: Your code here.
23         addl $8, %esp
24         popal
25
26         // Restore eflags from the stack.  After you do this, you can
27         // no longer use arithmetic operations or anything else that
28         // modifies eflags.
29         // LAB 4: Your code here.
30         addl $4, %esp
31         popfl
32
33         // Switch back to the adjusted trap-time stack.
34         // LAB 4: Your code here.
35         popl %esp
36
37         // Return to re-execute the instruction that faulted.
38         // LAB 4: Your code here.
39         ret
```

### 3.2.5.3  Exercise 11: Finish set_pgfault_handler() in lib/pgfault.c.

Firstly, see its hints:

```
1  // User-level page fault handler support.
2  // Rather than register the C page fault handler directly with the
3  // kernel as the page fault handler, we register the assembly language
4  // wrapper in pfentry.S, which in turns calls the registered C
5  // function.
6
7  #include <inc/lib.h>
8
9
10 // Assembly language pgfault entrypoint defined in lib/pfentry.S.
11 extern void _pgfault_upcall(void);
12
13 // Pointer to currently installed C-language pgfault handler.
14 void (*_pgfault_handler)(struct UTrapframe *utf);
15
16 //
17 // Set the page fault handler function.
18 // If there isn't one yet, _pgfault_handler will be 0.
```

```
19  // The first time we register a handler, we need to
20  // allocate an exception stack (one page of memory with its top
21  // at UXSTACKTOP), and tell the kernel to call the assembly-language
22  // _pgfault_upcall routine when a page fault occurs.
23  //
24  void
25  set_pgfault_handler(void (*handler)(struct UTrapframe *utf))
26  {
27          int r;
28
29          if (_pgfault_handler == 0) {
30                  // First time through!
31                  // LAB 4: Your code here.
32                  panic("set_pgfault_handler not implemented");
33          }
34
35          // Save handler pointer for assembly to call.
36          _pgfault_handler = handler;
37  }
```

See the implementation:

```
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                   if((r = sys_page_alloc(thisenv->env_id, (void *)(UXSTACKTOP - PGSIZE),
    ↪  PTE_P | PTE_W | PTE_U)) < 0)
10                          panic("set_pgfault_handler failed: %e", r);
11                  sys_env_set_pgfault_upcall(thisenv->env_id, _pgfault_upcall);
12                  //panic("set_pgfault_handler not implemented");
13          }
14
15          // Save handler pointer for assembly to call.
16          _pgfault_handler = handler;
17  }
```

Remind register the **sys_env_set_pgfault_upcall()** in **kern/syscall.c**:

```
1   // Set the page fault upcall for 'envid' by modifying the corresponding struct
2   // Env's 'env_pgfault_upcall' field.  When 'envid' causes a page fault, the
3   // kernel will push a fault record onto the exception stack, then branch to
4   // 'func'.
5   //
6   // Returns 0 on success, < 0 on error.  Errors are:
7   //      -E_BAD_ENV if environment envid doesn't currently exist,
8   //              or the caller doesn't have permission to change envid.
9   static int
10  sys_env_set_pgfault_upcall(envid_t envid, void *func)
11  {
12          // LAB 4: Your code here.
13          struct Env *env;
14          // the environment envid doesn't currently exist
15          // or the caller has no permission to change envid
16          if(envid2env(envid, &env, 1) < 0)
```

```
17            return -E_BAD_ENV;
18        // set the env_pgfault_upcall
19        env->env_pgfault_upcall = func;
20        return 0;
21
22        //panic("sys_env_set_pgfault_upcall not implemented");
23 }
```

### 3.2.6 Remember to Modify Your syscall() in kern/syscall.c

Add a case in **syscall()** in **kern/syscall.c**:

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

## 3.3 Implementing Copy-on-Write Fork

Like **dumbfork()**, **fork()** should:

- Create a new environment.

- Then scan through the parent environment's entire address space and set up corresponding page mappings in the child.

The key difference is that, while **dumbfork()** copied pages, **fork()** will initially only copy page mappings. **fork()** will copy each page only when one of the environments tries to write it.

The basic control flow for **fork()**:

The basic control flow for `fork()` is as follows:

1. The parent installs `pgfault()` as the C–level page fault handler, using the `set_pgfault_handler()` function you implemented above.
2. The parent calls `sys_exofork()` to create a child environment.
3. For each writable or copy–on–write page in its address space below UTOP, the parent calls `duppage`, which should map the page copy–on–write into the address space of the child and then *remap* the page copy–on–write in its own address space. [ Note: The ordering here (i.e., marking a page as COW in the child before marking it in the parent) actually matters! Can you see why? Try to think of a specific case where reversing the order could cause trouble. ] `duppage` sets both PTEs so that the page is not writeable, and to contain `PTE_COW` in the "avail" field to distinguish copy–on–write pages from genuine read–only pages.

   The exception stack is *not* remapped this way, however. Instead you need to allocate a fresh page in the child for the exception stack. Since the page fault handler will be doing the actual copying and the page fault handler runs on the exception stack, the exception stack cannot be made copy–on–write: who would copy it?

   `fork()` also needs to handle pages that are present, but not writable or copy–on–write.

4. The parent sets the user page fault entrypoint for the child to look like its own.
5. The child is now ready to run, so the parent marks it runnable.

Figure 13: Basic control flow of fork()

Each time one of the environments writes a copy-on-write page that it hasn't yet written, it will take a page fault. Here's the control flow for the user page fault handler:

The basic control flow for `fork()` is as follows:

1. The parent installs `pgfault()` as the C–level page fault handler, using the `set_pgfault_handler()` function you implemented above.
2. The parent calls `sys_exofork()` to create a child environment.
3. For each writable or copy–on–write page in its address space below UTOP, the parent calls `duppage`, which should map the page copy–on–write into the address space of the child and then *remap* the page copy–on–write in its own address space. [ Note: The ordering here (i.e., marking a page as COW in the child before marking it in the parent) actually matters! Can you see why? Try to think of a specific case where reversing the order could cause trouble. ] `duppage` sets both PTEs so that the page is not writeable, and to contain `PTE_COW` in the "avail" field to distinguish copy–on–write pages from genuine read–only pages.

   The exception stack is *not* remapped this way, however. Instead you need to allocate a fresh page in the child for the exception stack. Since the page fault handler will be doing the actual copying and the page fault handler runs on the exception stack, the exception stack cannot be made copy–on–write: who would copy it?

   `fork()` also needs to handle pages that are present, but not writable or copy–on–write.

4. The parent sets the user page fault entrypoint for the child to look like its own.
5. The child is now ready to run, so the parent marks it runnable.

Figure 14: Control flow of user page fault handler

### 3.3.1 Exercise 12

The user–level `lib/fork.c` code must consult the environment's page tables for several of the operations above (e.g., that the PTE for a page is marked `PTE_COW`). The kernel maps the environment's page tables at `UVPT` exactly for this purpose. It uses a [clever mapping trick](clever mapping trick) to make it to make it easy to lookup PTEs for user code. `lib/entry.s` sets up `uvpt` and `uvpd` so that you can easily lookup page–table information in `lib/fork.c`.

> **Exercise 12.** Implement `fork`, `duppage` and `pgfault` in `lib/fork.c`.
>
> Test your code with the `forktree` program. It should produce the following messages, with interspersed 'new env', 'free env', and 'exiting gracefully' messages. The messages may not appear in this order, and the environment IDs may be different.
>
> ```
> 1000: I am ''
> 1001: I am '0'
> 2000: I am '00'
> 2001: I am '000'
> 1002: I am '1'
> 3000: I am '11'
> 3001: I am '10'
> 4000: I am '100'
> 1003: I am '01'
> 5000: I am '010'
> 4001: I am '011'
> 2002: I am '110'
> 1004: I am '001'
> 1005: I am '111'
> 1006: I am '101'
> ```

Figure 15: Exercise 12

### 3.3.2 Clever

### 3.3.3 The Implementation of pgfault()

#### 3.3.3.1 Its Hints

```
// Custom page fault handler - if faulting page is copy-on-write,
// map in our own private writable copy.
//
static void
pgfault(struct UTrapframe *utf)
{
        void *addr = (void *) utf->utf_fault_va;
        uint32_t err = utf->utf_err;
        int r;

        // Check that the faulting access was (1) a write, and (2) to a
        // copy-on-write page.  If not, panic.
        // Hint:
        //   Use the read-only page table mappings at uvpt
        //   (see <inc/memlayout.h>).

        // LAB 4: Your code here.

        // Allocate a new page, map it at a temporary location (PFTEMP),
        // copy the data from the old page to the new page, then move the new
        // page to the old page's address.
        // Hint:
        //   You should make three system calls.

        // LAB 4: Your code here.

        panic("pgfault not implemented");
}
```

#### 3.3.3.2 Its Implementation

```
static void
pgfault(struct UTrapframe *utf)
{
        void *addr = (void *) utf->utf_fault_va;
        uint32_t err = utf->utf_err;
        int r;
        // LAB 4: Your code here.
        // Check that the faulting access was
        // (1). a write
        // (2). to a copy-on-write page.
        if((err & FEC_WR) == 0 || (uvpt[PGNUM(addr)] & PTE_COW) == 0)
                panic("pgfault failed: this page is not writable or attempt to access
    ↪ a non-cow page!");
        envid_t envid = sys_getenvid();
        if((r = sys_page_alloc(envid, (void *)PFTEMP, PTE_P|PTE_W|PTE_U)) < 0)
                panic("pgfault failed: page allocation failed!");

        addr = ROUNDDOWN(addr, PGSIZE);
        memmove(PFTEMP, addr, PGSIZE);
        if ((r = sys_page_unmap(envid, addr)) < 0)
                panic("pgfault failed: page unmap failed: %e", r);
```

```
21        if ((r = sys_page_map(envid, PFTEMP, envid, addr, PTE_P | PTE_W | PTE_U)) < 0)
22                panic("pgfault failed: page map failed: %e", r);
23        if ((r = sys_page_unmap(envid, PFTEMP)) < 0)
24                panic("pgfault failed: page map failed: %e", r);
25        //panic("pgfault not implemented");
26
27 }
```

### 3.3.4   The Implementation of duppage()

### 3.3.4.1   See its hints

```
1  // Map our virtual page pn (address pn*PGSIZE) into the target envid
2  // at the same virtual address.  If the page is writable or copy-on-write,
3  // the new mapping must be created copy-on-write, and then our mapping must be
4  // marked copy-on-write as well.  (Exercise: Why do we need to mark ours
5  // copy-on-write again if it was already copy-on-write at the beginning of
6  // this function?)
7  //
8  // Returns: 0 on success, < 0 on error.
9  // It is also OK to panic on error.
10 //
11 static int
12 duppage(envid_t envid, unsigned pn)
13 {
14        int r;
15
16        // LAB 4: Your code here.
17        panic("duppage not implemented");
18        return 0;
19 }
```

### 3.3.4.2   See its Implementation

```
1  static int
2  duppage(envid_t envid, unsigned pn)
3  {
4         int r;
5
6         // LAB 4: Your code here.
7         void *addr;
8         pte_t pte;
9         int perm;
10
11        addr = (void *)((uint32_t)pn * PGSIZE);
12        pte = uvpt[pn];
13        perm = PTE_P | PTE_U;
14        if ((pte & PTE_W) || (pte & PTE_COW))
15                perm |= PTE_COW;
16        if ((r = sys_page_map(thisenv->env_id, addr, envid, addr, perm)) < 0) {
17                panic("duppage failed: page remapping failed: %e", r);
18                return r;
19        }
20        if (perm & PTE_COW) {
21                if ((r = sys_page_map(thisenv->env_id, addr, thisenv->env_id, addr,
  ↪ perm)) < 0) {
22                        panic("duppage failed: page remapping failed: %e", r);
23                        return r;
```

47

```
24              }
25          }
26          //panic("duppage not implemented");
27          return 0;
28 }
```

### 3.3.5  User-Level Fork with Copy-on-Write Implementation

#### 3.3.5.1  See its Hints

```
1  //
2  // User-level fork with copy-on-write.
3  // Set up our page fault handler appropriately.
4  // Create a child.
5  // Copy our address space and page fault handler setup to the child.
6  // Then mark the child as runnable and return.
7  //
8  // Returns: child's envid to the parent, 0 to the child, < 0 on error.
9  // It is also OK to panic on error.
10 //
11 // Hint:
12 //   Use uvpd, uvpt, and duppage.
13 //   Remember to fix "thisenv" in the child process.
14 //   Neither user exception stack should ever be marked copy-on-write,
15 //   so you must allocate a new page for the child's user exception stack.
16 //
17 envid_t
18 fork(void)
19 {
20         // LAB 4: Your code here.
21         panic("fork not implemented");
22 }
```

#### 3.3.5.2  See its Implementation

```
1  envid_t
2  fork(void)
3  {
4          // LAB 4: Your code here.
5          envid_t childenvid;
6          uint32_t addr;
7          int i, j, pn, r;
8          extern void _pgfault_upcall(void);
9
10         // (1). Install pgfault() as the C-level page fault handler.
11         set_pgfault_handler(pgfault);
12
13         // (2). The parent calls sys_exofork() to create a child environment
14         if ((childenvid = sys_exofork()) < 0) {
15                 panic("sys_exofrok failed: %e", childenvid);
16                 return childenvid;
17         }
18         if(childenvid == 0) {
19                 // child environment
20                 thisenv = &envs[ENVX(sys_getenvid())];
21                 return 0;
22         }
```

```
23          // (3). For each writable or copy-on-write page in its address space before
   ↪ UTOP, the parent calls duppage, which should
24          //        (3.1) map the page copy-on-write into the address space of the child.
25          //        (3.2) then, remap the page copy-on-write into its own address space.
   ↪ Mind the order of (3.1) and (3.2)
26          // Following is parent environment operation.
27
28          //        (3.1) Copy address space and page fault handler setup
29          //              to the child.
30          for (pn = PGNUM(UTEXT); pn < PGNUM(USTACKTOP); pn++) {
31                  if((uvpd[pn >> 10] & PTE_P) && (uvpt[pn] & PTE_P)) {
32                          if ((r = duppage(childenvid, pn)) < 0)
33                                  return r;
34                  }
35          }
36
37          // allocate a page for child exception stack
38          if ((r = sys_page_alloc(childenvid, (void *)(UXSTACKTOP-PGSIZE), PTE_U | PTE_P
   ↪ | PTE_W)) < 0)
39                  return r;
40          // (4) The parent sets the user page fault entrypoint for
41          //      the child to look like its own.
42          if ((r = sys_env_set_pgfault_upcall(childenvid, _pgfault_upcall)) < 0)
43                  return r;
44          // (5). The child is now ready to run, so the parent marks it runnable
45          if ((r = sys_env_set_status(childenvid, ENV_RUNNABLE)) < 0)
46                  panic("fork failed: sys_env_set_status failed: %e", r);
47          return childenvid;
48
49          //panic("fork not implemented");
50 }
```

## 3.4 Grade

Now, we have grades:

Part A: 5/5. Part B: 50. Part C: 0/25.

# 4 Preemptive Multitasking and Inter-Process communication (IPC)

In this part, two objectives:

- Modify the kernel to preempt uncooperative environments.

- Allow environments to pass messages to each other explicitly.

## 4.1 Clock Interrupts and Preemption

In order to allow the kernel to **preempt** a running environment, forcefully retaking control of the CPU from this running environment, we must extend the JOS kernel to support external hardware interrupts from the clock hardware.

### 4.1.1 Interrupt Discipline

### 4.1.2 Exercise 13

**Exercise 13.** Modify `kern/trapentry.s` and `kern/trap.c` to initialize the appropriate entries in the IDT and provide handlers for IRQs 0 through 15. Then modify the code in `env_alloc()` in `kern/env.c` to ensure that user environments are always run with interrupts enabled.

Also uncomment the `sti` instruction in `sched_halt()` so that idle CPUs unmask interrupts.

The processor never pushes an error code when invoking a hardware interrupt handler. You might want to re-read section 9.2 of the 80386 Reference Manual, or section 5.8 of the IA-32 Intel Architecture Software Developer's Manual, Volume 3, at this time.

After doing this exercise, if you run your kernel with any test program that runs for a non-trivial length of time (e.g., `spin`), you should see the kernel print trap frames for hardware interrupts. While interrupts are now enabled in the processor, JOS isn't yet handling them, so you should see it misattribute each interrupt to the currently running user environment and destroy it. Eventually it should run out of environments to destroy and drop into the monitor.

Figure 16: Exercise 13

#### 4.1.2.1 Declaration and set SETGATE in trap_init() in kern/trap.c

```
// IRQs for LAB 4
void IRQ0();
SETGATE(idt[IRQ_OFFSET + 0], 0, GD_KT, IRQ0, 0);
```

```
4        void IRQ1();
5        SETGATE(idt[IRQ_OFFSET + 1], 0, GD_KT, IRQ1, 0);
6        void IRQ2();
7        SETGATE(idt[IRQ_OFFSET + 2], 0, GD_KT, IRQ2, 0);
8        void IRQ3();
9        SETGATE(idt[IRQ_OFFSET + 3], 0, GD_KT, IRQ3, 0);
10       void IRQ4();
11       SETGATE(idt[IRQ_OFFSET + 4], 0, GD_KT, IRQ4, 0);
12       void IRQ5();
13       SETGATE(idt[IRQ_OFFSET + 5], 0, GD_KT, IRQ5, 0);
14       void IRQ6();
15       SETGATE(idt[IRQ_OFFSET + 6], 0, GD_KT, IRQ6, 0);
16       void IRQ7();
17       SETGATE(idt[IRQ_OFFSET + 7], 0, GD_KT, IRQ7, 0);
18       void IRQ8();
19       SETGATE(idt[IRQ_OFFSET + 8], 0, GD_KT, IRQ8, 0);
20       void IRQ9();
21       SETGATE(idt[IRQ_OFFSET + 9], 0, GD_KT, IRQ9, 0);
22       void IRQ10();
23       SETGATE(idt[IRQ_OFFSET + 10], 0, GD_KT, IRQ10, 0);
24       void IRQ11();
25       SETGATE(idt[IRQ_OFFSET + 11], 0, GD_KT, IRQ11, 0);
26       void IRQ12();
27       SETGATE(idt[IRQ_OFFSET + 12], 0, GD_KT, IRQ12, 0);
28       void IRQ13();
29       SETGATE(idt[IRQ_OFFSET + 13], 0, GD_KT, IRQ13, 0);
30       void IRQ14();
31       SETGATE(idt[IRQ_OFFSET + 14], 0, GD_KT, IRQ14, 0);
32       void IRQ15();
33       SETGATE(idt[IRQ_OFFSET + 15], 0, GD_KT, IRQ15, 0);
```

#### 4.1.2.2 Generating entry points for IRQs in kern/trapentry.S

```
1  // LAB 4
2  // IRQs
3  // LAB 4
4  // IRQs
5  TRAPHANDLER_NOEC(IRQ0, IRQ_OFFSET+0);
6  TRAPHANDLER_NOEC(IRQ1, IRQ_OFFSET+1);
7  TRAPHANDLER_NOEC(IRQ2, IRQ_OFFSET+2);
8  TRAPHANDLER_NOEC(IRQ3, IRQ_OFFSET+3);
9  TRAPHANDLER_NOEC(IRQ4, IRQ_OFFSET+4);
10 TRAPHANDLER_NOEC(IRQ5, IRQ_OFFSET+5);
11 TRAPHANDLER_NOEC(IRQ6, IRQ_OFFSET+6);
12 TRAPHANDLER_NOEC(IRQ7, IRQ_OFFSET+7);
13 TRAPHANDLER_NOEC(IRQ8, IRQ_OFFSET+8);
14 TRAPHANDLER_NOEC(IRQ9, IRQ_OFFSET+9);
15 TRAPHANDLER_NOEC(IRQ10, IRQ_OFFSET+10);
16 TRAPHANDLER_NOEC(IRQ11, IRQ_OFFSET+11);
17 TRAPHANDLER_NOEC(IRQ12, IRQ_OFFSET+12);
18 TRAPHANDLER_NOEC(IRQ13, IRQ_OFFSET+13);
19 TRAPHANDLER_NOEC(IRQ14, IRQ_OFFSET+14);
20 TRAPHANDLER_NOEC(IRQ15, IRQ_OFFSET+15);
```

#### 4.1.2.3 Modify env_alloc() in kern/env.c to ensure that user environment are always run with interrupts enabled.

```
1  // Allocates and initializes a new environment.
```

51

```c
// On success, the new environment is stored in *newenv_store.
//
// Returns 0 on success, < 0 on failure.  Errors include:
//      -E_NO_FREE_ENV if all NENV environments are allocated
//      -E_NO_MEM on memory exhaustion
//
int
env_alloc(struct Env **newenv_store, envid_t parent_id)
{
        int32_t generation;
        int r;
        struct Env *e;

        if (!(e = env_free_list))
                return -E_NO_FREE_ENV;

        // Allocate and set up the page directory for this environment.
        if ((r = env_setup_vm(e)) < 0)
                return r;

        // Generate an env_id for this environment.
        generation = (e->env_id + (1 << ENVGENSHIFT)) & ~(NENV - 1);
        if (generation <= 0)    // Don't create a negative env_id.
                generation = 1 << ENVGENSHIFT;
        e->env_id = generation | (e - envs);

        // Set the basic status variables.
        e->env_parent_id = parent_id;
        e->env_type = ENV_TYPE_USER;
        e->env_status = ENV_RUNNABLE;
        e->env_runs = 0;
        // Clear out all the saved register state,
        // to prevent the register values
        // of a prior environment inhabiting this Env structure
        // from "leaking" into our new environment.
        memset(&e->env_tf, 0, sizeof(e->env_tf));

        // Set up appropriate initial values for the segment registers.
        // GD_UD is the user data segment selector in the GDT, and
        // GD_UT is the user text segment selector (see inc/memlayout.h).
        // The low 2 bits of each segment register contains the
        // Requestor Privilege Level (RPL); 3 means user mode.  When
        // we switch privilege levels, the hardware does various
        // checks involving the RPL and the Descriptor Privilege Level
        // (DPL) stored in the descriptors themselves.
        e->env_tf.tf_ds = GD_UD | 3;
        e->env_tf.tf_es = GD_UD | 3;
        e->env_tf.tf_ss = GD_UD | 3;
        e->env_tf.tf_esp = USTACKTOP;
        e->env_tf.tf_cs = GD_UT | 3;
        // You will set e->env_tf.tf_eip later.

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

        // Clear the page fault handler until user installs one.
        e->env_pgfault_upcall = 0;
```

52

```
60
61          // Also clear the IPC receiving flag.
62          e->env_ipc_recving = 0;
63
64          // commit the allocation
65          env_free_list = e->env_link;
66          *newenv_store = e;
67
68          cprintf("[%08x] new env %08x\n", curenv ? curenv->env_id : 0, e->env_id);
69          return 0;
70 }
```

Just add one line:

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

### 4.1.2.4   DO NOT FORGET COMMENT the sti in sched_halt() in kern/sched.c

```
1  // Halt this CPU when there is nothing to do. Wait until the
2  // timer interrupt wakes it up. This function never returns.
3  //
4  void
5  sched_halt(void)
6  {
7          int i;
8
9          // For debugging and testing purposes, if there are no runnable
10         // environments in the system, then drop into the kernel monitor.
11         for (i = 0; i < NENV; i++) {
12                 if ((envs[i].env_status == ENV_RUNNABLE ||
13                     envs[i].env_status == ENV_RUNNING ||
14                     envs[i].env_status == ENV_DYING))
15                         break;
16         }
17         if (i == NENV) {
18                 cprintf("No runnable environments in the system!\n");
19                 while (1)
20                         monitor(NULL);
21         }
22
23         // Mark that no environment is running on this CPU
24         curenv = NULL;
25         lcr3(PADDR(kern_pgdir));
26
27         // Mark that this CPU is in the HALT state, so that when
28         // timer interupts come in, we know we should re-acquire the
29         // big kernel lock
30         // big kernel lock
31         xchg(&thiscpu->cpu_status, CPU_HALTED);
32
33         // Release the big kernel lock as if we were "leaving" the kernel
34         unlock_kernel();
35
36         // Reset stack pointer, enable interrupts and then halt.
37         asm volatile (
38                 "movl $0, %%ebp\n"
```

```
39            "movl %0, %%esp\n"
40            "pushl $0\n"
41            "pushl $0\n"
42            // Uncomment the following line after completing exercise 13
43            //"sti\n"
44            "1:\n"
45            "hlt\n"
46            "jmp 1b\n"
47        : : "a" (thiscpu->cpu_ts.ts_esp0));
48 }
```

### 4.1.3 Handling Clock Interrupts

### 4.1.4 Exercise 14

**Exercise 14.** Modify the kernel's `trap_dispatch()` function so that it calls `sched_yield()` to find and run a different environment whenever a clock interrupt takes place.

You should now be able to get the `user/spin` test to work: the parent environment should fork off the child, `sys_yield()` to it a couple times but in each case regain control of the CPU after one time slice, and finally kill the child environment and terminate gracefully.

Figure 17: Exercise 14

```
1 static void
2 trap_dispatch(struct Trapframe *tf)
3 {
4        // Handle spurious interrupts
5        // The hardware sometimes raises these because of noise on the
6        // IRQ line or other reasons. We don't care.
7        if (tf->tf_trapno == IRQ_OFFSET + IRQ_SPURIOUS) {
8                cprintf("Spurious interrupt on irq 7\n");
9                print_trapframe(tf);
10               return;
11        }
12        // Handle processor exceptions.
13        // LAB 3: Your code here.
14        switch(tf->tf_trapno) {
15                case T_PGFLT:
16                        page_fault_handler(tf);
17                        return;
18                case T_BRKPT:
19                        monitor(tf);
20                        return;
```

```
21              case T_SYSCALL:
22                      tf->tf_regs.reg_eax = syscall(tf->tf_regs.reg_eax,
23                              tf->tf_regs.reg_edx,
24                              tf->tf_regs.reg_ecx, tf->tf_regs.reg_ebx,
25                              tf->tf_regs.reg_edi, tf->tf_regs.reg_esi);
26                      return;
27              // Handle clock interrupts. Don't forget to acknowledge the
28              // interrupt using lapic_eoi() before calling the scheduler!
29              // LAB 4: Your code here.
30              case IRQ_OFFSET + IRQ_TIMER:
31                      lapic_eoi();
32                      sched_yield();
33                      return;
34              default:
35                      // Unexpected trap: The user process or the kernel has a bug.
36                      print_trapframe(tf);
37                      if (tf->tf_cs == GD_KT)
38                              panic("unhandled trap in kernel");
39                      else {
40                              env_destroy(curenv);
41                              return;
42                      }
43                      break;
44      }
45 }
```

### 4.1.5 Present Score

Now:

| Part   | Score  |
|--------|--------|
| Part A | 5/5    |
| Part B | 50/50  |
| Part C | 10/25  |
| Final  | 65/80  |

## 4.2 Inter-Process communication (IPC)

### 4.2.1 Background Knowledge about IPC

- We've been focusing on the isolation aspects of the operating system, the ways it provides the illusion that each program has a machine all to itself;

- Another important service of an operating system is to allow programs to communicate with each other when they want to.

The Unix pipe model is the canonical example.

You will implement two system calls, **sys_ipc_recv** and **sys_ipc_try_send**. Then you will implement two library wrappers **ipc_recv** and **ipc_send**.

The *message* that user environments can send consists of two components:

- A single 32-bit value.

- Optionally a single page mapping. Allowing environments to pass page mappings in messages provides an efficient way to transfer more data than will fit into a single 32-bit integer, and also allows environments to set up shared memory arrangements easily.

### 4.2.2 Exercise 15

**Exercise 15.** Implement `sys_ipc_recv` and `sys_ipc_try_send` in `kern/syscall.c`. Read the comments on both before implementing them, since they have to work together. When you call `envid2env` in these routines, you should set the `checkperm` flag to 0, meaning that any environment is allowed to send IPC messages to any other environment, and the kernel does no special permission checking other than verifying that the target envid is valid.

Then implement the `ipc_recv` and `ipc_send` functions in `lib/ipc.c`.

Use the `user/pingpong` and `user/primes` functions to test your IPC mechanism. `user/primes` will generate for each prime number a new environment until JOS runs out of environments. You might find it interesting to read `user/primes.c` to see all the forking and IPC going on behind the scenes.

Figure 18: Exercise 15

### 4.2.3 The Hints in kern/syscall.c

```
// Try to send 'value' to the target env 'envid'.
// If srcva < UTOP, then also send page currently mapped at 'srcva',
// so that receiver gets a duplicate mapping of the same page.
//
// The send fails with a return value of -E_IPC_NOT_RECV if the
// target is not blocked, waiting for an IPC.
//
// The send also can fail for the other reasons listed below.
//
// Otherwise, the send succeeds, and the target's ipc fields are
// updated as follows:
//    env_ipc_recving is set to 0 to block future sends;
//    env_ipc_from is set to the sending envid;
//    env_ipc_value is set to the 'value' parameter;
//    env_ipc_perm is set to 'perm' if a page was transferred, 0 otherwise.
// The target environment is marked runnable again, returning 0
// from the paused sys_ipc_recv system call.  (Hint: does the
// sys_ipc_recv function ever actually return?)
//
// If the sender wants to send a page but the receiver isn't asking for one,
// then no page mapping is transferred, but no error occurs.
// The ipc only happens when no errors occur.
//
// Returns 0 on success, < 0 on error.
```

56

```
25 // Errors are:
26 //         -E_BAD_ENV if environment envid doesn't currently exist.
27 //                 (No need to check permissions.)
28 //         -E_IPC_NOT_RECV if envid is not currently blocked in sys_ipc_recv,
29 //                 or another environment managed to send first.
30 //         -E_INVAL if srcva < UTOP but srcva is not page-aligned.
31 //         -E_INVAL if srcva < UTOP and perm is inappropriate
32 //                 (see sys_page_alloc).
33 //         -E_INVAL if srcva < UTOP but srcva is not mapped in the caller's
34 //                 address space.
35 //         -E_INVAL if (perm & PTE_W), but srcva is read-only in the
36 //                 current environment's address space.
37 //         -E_NO_MEM if there's not enough memory to map srcva in envid's
38 //                 address space.
39 static int
40 sys_ipc_try_send(envid_t envid, uint32_t value, void *srcva, unsigned perm)
41 {
42         // LAB 4: Your code here.
43         panic("sys_ipc_try_send not implemented");
44 }
45
46 // Block until a value is ready.  Record that you want to receive
47 // using the env_ipc_recving and env_ipc_dstva fields of struct Env,
48 // mark yourself not runnable, and then give up the CPU.
49 //
50 // If 'dstva' is < UTOP, then you are willing to receive a page of data.
51 // 'dstva' is the virtual address at which the sent page should be mapped.
52 //
53 // This function only returns on error, but the system call will eventually
54 // return 0 on success.
55 // Return < 0 on error.  Errors are:
56 //         -E_INVAL if dstva < UTOP but dstva is not page-aligned.
57 static int
58 sys_ipc_recv(void *dstva)
59 {
60         // LAB 4: Your code here.
61         panic("sys_ipc_recv not implemented");
62         return 0;
63 }
```

## 4.2.4 Sending and Receiving Message

### 4.2.4.1 To receive a message

To receive a message, an environment calls **sys_ipc_recv**. This system call **de-schedules** the current environment and does not run it again until a message has been received.

Here is a safe technique for **malfunction**:

When an environment is waiting to receive a message, **any other** environment can send it a message. In this *safe* approach, an environment cannot cause another environment to malfunction simply by sending a message.

### 4.2.4.2 To try to send a value

To try to send a value, an environment calls **sys_ipc_try_send** with the {the receiver environment id, and the value to be sent}. Then they are two cases in this situation:

- If the named environment is actually receiving, then the sender delivers the message and return 0.

- Otherwise, the send returns **-E_IPC_NOT_RECV** to indicate the target environment is not currently expecting to receive a value.

### 4.2.5 Transferring Pages

#### 4.2.5.1 Receive a page mapping

When an environment calls **sys_ipc_recv** with a valid **dstva**(below **UTOP**) parameter, the environment is stating that it is willing to receive a page mapping. If the sender sends a page, then that page should be mapped at the **dstva** in the receiver's address space. If the receiver already had a page mapped at **dstva**, then that previous page is unmapped.

#### 4.2.5.2 Send a page mapping

When an environment calls **sys_ipc_try_send** with a valid **srcva**, it means that the sender wants to send the page currently mapped at **srcva** to the receiver, with permission **perm**. After a successful IPC:

- The sender keeps its original mapping for the page at **srcva** in its address space.

- The receiver also obtains a mapping for this same physical page at the **dstva** specified by the receiver in the receiver's address space.

- This physical page is shared between the sender and the receiver.

If either the sender or the receiver does not indicate that a page should be transferred, then no page is transferred. After any IPC the kernel sets the new field **env_ipc_perm** in the receiver's **Env** structure to the permissions of the page received or 0 if no page was received.

### 4.2.6 The Implementation of sys_ipc_try_send() in kern/syscall.c

#### 4.2.6.1 Let's see its hints

```
// Try to send 'value' to the target env 'envid'.
// If srcva < UTOP, then also send page currently mapped at 'srcva',
// so that receiver gets a duplicate mapping of the same page.
//
// The send fails with a return value of -E_IPC_NOT_RECV if the
// target is not blocked, waiting for an IPC.
//
// The send also can fail for the other reasons listed below.
//
// Otherwise, the send succeeds, and the target's ipc fields are
// updated as follows:
//    env_ipc_recving is set to 0 to block future sends;
//    env_ipc_from is set to the sending envid;
//    env_ipc_value is set to the 'value' parameter;
//    env_ipc_perm is set to 'perm' if a page was transferred, 0 otherwise.
// The target environment is marked runnable again, returning 0
// from the paused sys_ipc_recv system call.  (Hint: does the
```

```
18 // sys_ipc_recv function ever actually return?)
19 //
20 // If the sender wants to send a page but the receiver isn't asking for one,
21 // then no page mapping is transferred, but no error occurs.
22 // The ipc only happens when no errors occur.
23 //
24 // Returns 0 on success, < 0 on error.
25 // Errors are:
26 //       -E_BAD_ENV if environment envid doesn't currently exist.
27 //             (No need to check permissions.)
28 //       -E_IPC_NOT_RECV if envid is not currently blocked in sys_ipc_recv,
29 //             or another environment managed to send first.
30 //       -E_INVAL if srcva < UTOP but srcva is not page-aligned.
31 //       -E_INVAL if srcva < UTOP and perm is inappropriate
32 //             (see sys_page_alloc).
33 //       -E_INVAL if srcva < UTOP but srcva is not mapped in the caller's
34 //             address space.
35 //       -E_INVAL if (perm & PTE_W), but srcva is read-only in the
36 //             current environment's address space.
37 //       -E_NO_MEM if there's not enough memory to map srcva in envid's
38 //             address space.
39 static int
40 sys_ipc_try_send(envid_t envid, uint32_t value, void *srcva, unsigned perm)
41 {
42         // LAB 4: Your code here.
43         panic("sys_ipc_try_send not implemented");
44 }
```

### 4.2.6.2   See the Implementation

```
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         struct Env *dstenv;
6         pte_t *pte;
7         struct PageInfo *pp;
8         int r;
9         // environment envid doesn't currently exist.
10        if ((r = envid2env(envid, &dstenv, 0)) < 0)
11                return -E_BAD_ENV;
12        //if envid is not currently blocked in sys_ipc_recv
13        if((dstenv->env_ipc_recving != true) || (dstenv->env_ipc_from != 0))
14                return -E_IPC_NOT_RECV;
15        // if srcva < UTOP but srcva is not page-aligned
16        if ((srcva < (void *)UTOP) && (PGOFF(srcva)))
17                return -E_INVAL;
18        // if srcva < UTOP but perm is inappropriate
19        if (srcva < (void *)UTOP) {
20                // perm inappropriate
21                if(((perm & PTE_P) == 0) || ((perm & PTE_U) == 0))
22                        return -E_INVAL;
23                if ((perm & ~(PTE_P | PTE_U | PTE_W | PTE_AVAIL)) != 0)
24                        return -E_INVAL;
25                // no page mapped at srcva
26                if ((pp = page_lookup(curenv->env_pgdir, srcva, &pte)) == NULL)
27                        return -E_INVAL;
28                // (perm & PTE_W), but srcva is read-only
```

59

```
29          if(((perm & PTE_W) != 0) && ((*pte & PTE_W) == 0))
30                  return -E_INVAL;
31          if (dstenv->env_ipc_dstva != 0) {
32                  // there is no enough memory to map srcva
33                  // in envid's address space.
34                  if ((r = page_insert(dstenv->env_pgdir, pp, dstenv->
   ↪ env_ipc_dstva, perm)) < 0)
35                          return -E_NO_MEM;
36                  dstenv->env_ipc_perm = perm;
37          }
38
39      }
40      dstenv->env_ipc_from = curenv->env_id;
41      dstenv->env_ipc_recving = false;
42      dstenv->env_ipc_value = value;
43      dstenv->env_status = ENV_RUNNABLE;
44      dstenv->env_tf.tf_regs.reg_eax = 0;
45      return 0;
46
47      //panic("sys_ipc_try_send not implemented");
48  }
```

### 4.2.7   The Implementation of sys_ipc_recv() in kern/syscall.c

#### 4.2.7.1   See its hints

```
1  // Block until a value is ready.  Record that you want to receive
2  // using the env_ipc_recving and env_ipc_dstva fields of struct Env,
3  // mark yourself not runnable, and then give up the CPU.
4  //
5  // If 'dstva' is < UTOP, then you are willing to receive a page of data.
6  // 'dstva' is the virtual address at which the sent page should be mapped.
7  //
8  // This function only returns on error, but the system call will eventually
9  // return 0 on success.
10 // Return < 0 on error.  Errors are:
11 //      -E_INVAL if dstva < UTOP but dstva is not page-aligned.
12 static int
13 sys_ipc_recv(void *dstva)
14 {
15     // LAB 4: Your code here.
16     panic("sys_ipc_recv not implemented");
17     return 0;
18 }
```

#### 4.2.7.2   See the implementation

```
1  static int
2  sys_ipc_recv(void *dstva)
3  {
4      // LAB 4: Your code here.
5      if (dstva < (void *)UTOP && PGOFF(dstva))
6              return -E_INVAL;
7      // Record that you want to receive using env_ipc_recving
8      // and env_ipc_dstva fields of struct Env
9      curenv->env_ipc_recving = true;
10     curenv->env_ipc_dstva = dstva;
11     // mark yourself not runnable.
```

```
12          curenv->env_status = ENV_NOT_RUNNABLE;
13          curenv->env_ipc_from = 0;
14          // Give up CPU
15          sched_yield();
16          return 0;
17          //panic("sys_ipc_recv not implemented");
18   }
```

### 4.2.8   DO NOT FORGET REGISTER IN syscall() in kern/syscall.c

```
1          // LAB 4: sys_ipc_recv and sys_ipc_try_send
2          case SYS_ipc_recv:
3                  return sys_ipc_recv((void *)a1);
4          case SYS_ipc_try_send:
5                  return sys_ipc_try_send((envid_t)a1, (uint32_t)a2, (void *)a3, (
    ↪ unsigned)a4);
```

### 4.2.9   The Implementation of ipc_recv() in kern/ipc.c

#### 4.2.9.1   Firstly, see its hints

```
1  // Receive a value via IPC and return it.
2  // If 'pg' is nonnull, then any page sent by the sender will be mapped at
3  //      that address.
4  // If 'from_env_store' is nonnull, then store the IPC sender's envid in
5  //      *from_env_store.
6  // If 'perm_store' is nonnull, then store the IPC sender's page permission
7  //      in *perm_store (this is nonzero iff a page was successfully
8  //      transferred to 'pg').
9  // If the system call fails, then store 0 in *fromenv and *perm (if
10 //      they're nonnull) and return the error.
11 // Otherwise, return the value sent by the sender
12 //
13 // Hint:
14 //   Use 'thisenv' to discover the value and who sent it.
15 //   If 'pg' is null, pass sys_ipc_recv a value that it will understand
16 //   as meaning "no page".  (Zero is not the right value, since that's
17 //   a perfectly valid place to map a page.)
18 int32_t
19 ipc_recv(envid_t *from_env_store, void *pg, int *perm_store)
20 {
21          // LAB 4: Your code here.
22          panic("ipc_recv not implemented");
23          return 0;
24 }
```

#### 4.2.9.2   See its implementation

```
1  int32_t
2  ipc_recv(envid_t *from_env_store, void *pg, int *perm_store)
3  {
4          // LAB 4: Your code here.
5          int r;
6          if (pg == NULL)
7                  pg = (void *)UTOP;
8          r = sys_ipc_recv(pg);
9          // If 'from_env_store' is nonnull, then store the IPC
10         // sender's envid in *from_env_store
```

```
11          if (from_env_store != NULL)
12                  *from_env_store = r < 0 ? 0 : thisenv->env_ipc_
13 from;
14          // If 'perm store' is nonnull, then store the IPC sende
15 r's
16          // page permission in *perm_store
17          if (perm_store != NULL)
18                  *perm_store = r < 0 ? 0 : thisenv->env_ipc_perm
19 ;
20
21          // If the system call fails, store 0 in *fromenv and
22          // *perm if they're nonnull and return the error.
23          // Otherwise, return the value sent by the sender
24          if (r < 0)
25                  return r;
26          else
27                  return thisenv->env_ipc_value;
28          //panic("ipc_recv not implemented");
29          //return 0;
30 }
```

## 4.2.10 The Implementation of ipc_send() in kern/ipc.c

### 4.2.10.1 Hints of ipc_send()

```
1  // Send 'val' (and 'pg' with 'perm', if 'pg' is nonnull) to 'toenv'.
2  // This function keeps trying until it succeeds.
3  // It should panic() on any error other than -E_IPC_NOT_RECV.
4  //
5  // Hint:
6  //   Use sys_yield() to be CPU-friendly.
7  //   If 'pg' is null, pass sys_ipc_try_send a value that it will understand
8  //   as meaning "no page".  (Zero is not the right value.)
9  void
10 ipc_send(envid_t to_env, uint32_t val, void *pg, int perm)
11 {
12         // LAB 4: Your code here.
13         panic("ipc_send not implemented");
14 }
```

### 4.2.10.2 See its implementation

```
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          void *dstpg;
7
8          dstpg = (pg == NULL ? (void *)UTOP : pg);
9          // Keep trying unitl it succeeds.
10         while ((r = sys_ipc_try_send(to_env, val, dstpg, perm)) < 0) {
11                 if (r != -E_IPC_NOT_RECV)
12                         panic("ipc_send failed: error other than -E_IPC_NOT_RECV
    ↪ happened: %e", r);
13                 // Use sys_yield() to be CPU-friendly
14                 sys_yield();
15         }
```

```
16          //panic("ipc_send not implemented");
17 }
```

## 4.3  Run user/pingpong and user/primes to Test IPC

This result is useful for check the merge result of LAB 4 and LAB 5. Just see the result:

### 4.3.1  Run user/pingpong

#### 4.3.1.1  See the code in pingpong.c

```
1 // Ping-pong a counter between two processes.
2 // Only need to start one of these -- splits into two with fork.
3
4 #include <inc/lib.h>
5
6 void
7 umain(int argc, char **argv)
8 {
9          envid_t who;
10
11         if ((who = fork()) != 0) {
12                 // get the ball rolling
13                 cprintf("send 0 from %x to %x\n", sys_getenvid(), who);
14                 ipc_send(who, 0, 0, 0);
15         }
16
17         while (1) {
18                 uint32_t i = ipc_recv(&who, 0, 0);
19                 cprintf("%x got %d from %x\n", sys_getenvid(), i, who);
20                 if (i == 10)
21                         return;
22                 i++;
23                 ipc_send(who, i, 0, 0);
24                 if (i == 10)
25                         return;
26         }
27 }
```

#### 4.3.1.2  See the results

```
1 make[1]: Entering directory '/home/cui/mit6828/lab'
2 + cc kern/init.c
3 + cc[USER] user/pingpong.c
4 + ld obj/user/pingpong
5 + ld obj/kern/kernel
6 + mk obj/kern/kernel.img
7 make[1]: Leaving directory '/home/cui/mit6828/lab'
8 qemu-system-i386 -drive file=obj/kern/kernel.img,index=0,media=disk,format=raw -serial
    ↪   mon:stdio -gdb tcp::26000 -D qemu.log -smp 1
9 6828 decimal is 15254 octal!
10 Physical memory: 131072K available, base = 640K, extended = 130432K
11 check_page_free_list() succeeded!
12 check_page_alloc() succeeded!
13 check_page() succeeded!
14 check_kern_pgdir() succeeded!
15 check_page_free_list() succeeded!
```

```
16  check_page_installed_pgdir() succeeded!
17  SMP: CPU 0 found 1 CPU(s)
18  enabled interrupts: 1 2
19  [00000000] new env 00001000
20  [00001000] new env 00001001
21  send 0 from 1000 to 1001
22  1001 got 0 from 1000
23  1000 got 1 from 1001
24  1001 got 2 from 1000
25  1000 got 3 from 1001
26  1001 got 4 from 1000
27  1000 got 5 from 1001
28  1001 got 6 from 1000
29  1000 got 7 from 1001
30  1001 got 8 from 1000
31  1000 got 9 from 1001
32  [00001000] exiting gracefully
33  [00001000] free env 00001000
34  1001 got 10 from 1000
35  [00001001] exiting gracefully
36  [00001001] free env 00001001
37  No runnable environments in the system!
38  Welcome to the JOS kernel monitor!
39  Type 'help' for a list of commands.
40  K>
```

### 4.3.2   Run user/primes

```
1   // Concurrent version of prime sieve of Eratosthenes.
2   // Invented by Doug McIlroy, inventor of Unix pipes.
3   // See http://swtch.com/~rsc/thread/.
4   // The picture halfway down the page and the text surrounding it
5   // explain what's going on here.
6   //
7   // Since NENV is 1024, we can print 1022 primes before running out.
8   // The remaining two environments are the integer generator at the bottom
9   // of main and user/idle.
10  #include <inc/lib.h>
11
12  unsigned
13  primeproc(void)
14  {
15          int i, id, p;
16          envid_t envid;
17
18          // fetch a prime from our left neighbor
19  top:
20          p = ipc_recv(&envid, 0, 0);
21          cprintf("CPU %d: %d ", thisenv->env_cpunum, p);
22
23          // fork a right neighbor to continue the chain
24          if ((id = fork()) < 0)
25                  panic("fork: %e", id);
26          if (id == 0)
27                  goto top;
28
29          // filter out multiples of our prime
30          while (1) {
```

64

```
31                  i = ipc_recv(&envid, 0, 0);
32                  if (i % p)
33                          ipc_send(id, i, 0, 0);
34          }
35  }
36
37  void
38  umain(int argc, char **argv)
39  {
40          int i, id;
41
42          // fork the first prime process in the chain
43          if ((id = fork()) < 0)
44                  panic("fork: %e", id);
45          if (id == 0)
46                  primeproc();
47
48          // feed all the integers through
49          for (i = 2; ; i++)
50                  ipc_send(id, i, 0, 0);
51  }
```

#### 4.3.2.1   The results

```
1   make[1]: Entering directory '/home/cui/mit6828/lab'
2   make[1]: Leaving directory '/home/cui/mit6828/lab'
3   qemu-system-i386 -drive file=obj/kern/kernel.img,index=0,media=disk,format=raw -serial
        ↪ mon:stdio -gdb tcp::26000 -D qemu.log -smp 1
4   6828 decimal is 15254 octal!
5   Physical memory: 131072K available, base = 640K, extended = 130432K
6   check_page_free_list() succeeded!
7   check_page_alloc() succeeded!
8   check_page() succeeded!
9   check_kern_pgdir() succeeded!
10  check_page_free_list() succeeded!
11  check_page_installed_pgdir() succeeded!
12  SMP: CPU 0 found 1 CPU(s)
13  enabled interrupts: 1 2
14  [00000000] new env 00001000
15  [00001000] new env 00001001
16  CPU 0: 2 [00001001] new env 00001002
17  CPU 0: 3 [00001002] new env 00001003
18  CPU 0: 5 [00001003] new env 00001004
19  CPU 0: 7 [00001004] new env 00001005
20  CPU 0: 11 [00001005] new env 00001006
21  CPU 0: 13 [00001006] new env 00001007
22  CPU 0: 17 [00001007] new env 00001008
23  CPU 0: 19 [00001008] new env 00001009
24  CPU 0: 23 [00001009] new env 0000100a
25  CPU 0: 29 [0000100a] new env 0000100b
26  CPU 0: 31 [0000100b] new env 0000100c
27  CPU 0: 37 [0000100c] new env 0000100d
28  CPU 0: 41 [0000100d] new env 0000100e
29  CPU 0: 43 [0000100e] new env 0000100f
30  CPU 0: 47 [0000100f] new env 00001010
31  CPU 0: 53 [00001010] new env 00001011
32  CPU 0: 59 [00001011] new env 00001012
33  CPU 0: 61 [00001012] new env 00001013
```

```
34 .........................................................
35 .........................................................
36 CPU 0: 8093 [000013fa] new env 000013fb
37 CPU 0: 8101 [000013fb] new env 000013fc
38 CPU 0: 8111 [000013fc] new env 000013fd
39 CPU 0: 8117 [000013fd] new env 000013fe
40 CPU 0: 8123 [000013fe] new env 000013ff
41 CPU 0: 8147 [000013ff] user panic in <unknown> at lib/fork.c:127: sys_exofrok failed:
     ↪ out of environments
42 Welcome to the JOS kernel monitor!
43 Type 'help' for a list of commands.
44 TRAP frame at 0xf02d3f84 from CPU 0
45   edi  0x00000000
46   esi  0x008016a3
47   ebp  0xeebfdf30
48   oesp 0xeffffffdc
49   ebx  0xeebfdf44
50   edx  0xeebfdddd8
51   ecx  0x00000001
52   eax  0x00000001
53   es   0x----0023
54   ds   0x----0023
55   trap 0x00000003 Breakpoint
56   err  0x00000000
57   eip  0x008001f8
58   cs   0x----001b
59   flag 0x00000292
60   esp  0xeebfdf08
61   ss   0x----0023
62 K>
```

## 4.4   Final Score

| Part   | Score |
|--------|-------|
| Part A | 5/5   |
| Part B | 50/50 |
| Part C | 25/25 |
| Final  | 80/80 |

# 5 Reference

1. bysui's github and blog

   - github: https://github.com/bysui/mit6.828
   - blog: https://blog.csdn.net/bysui

2. SmallPond's github and blog

   - github: https://github.com/SmallPond/MIT6.828_OS
   - blog: https://me.csdn.net/Small_Pond

3. SimpCosm's github

   - github: https://github.com/SimpCosm/6.828

4. fatsheep9146's blog

   - blog: https://www.cnblogs.com/fatsheep9146/category/769143.html