# MIT 6.828: Operating System Engineering

Lab Report
Lab 3: User Environments

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

# Contents

# 1   Introduction

## 1.1   Useful Flags of PTE

Each PTE contains flag bits that tell the paging hardware how the associated virtual address is allowed to be used.

1. **PTE_P**: indicates whether the PTE is present.

2. **PTE_W**: controls whether instructions are allowed to issue writes to the page.

3. **PTE_U**: controls whether user programs are allowed to use page.

Hardware wants attention, so software must set aside current work and respond. Why does hardware want attention now?

- MMU cannot translate address.

- User program divides by zero.

- User program wants to execute kernel code (INT).

- Kernel CPU-to-CPU communication, e.g., to flush TLB.

These traps, generally speaking, fall into 3 classes:

1. Exceptions (page fault, divide by zero).

2. System calls (INT, intended exception).

3. Interrupts (devices want attention).

**How does trap() know which device interrupted?**

1. Kernel tells LAPIC/IOAPIC what vector number to use.

2. IDT associates an instruction address with each vector number.

3. Execute the associated instruction.

4. Each vector jumps to alltraps.

# 2 XV6 Chapter 3: Traps, interrupts, and drivers

# 3 Part A: User Environments and Exception Handling

The kernel maintains three main global variables pertaining to environments:

```
struct Env *envs = NULL;              // All environments
struct Env *curenv = NULL;            // The current env
static struct Env *env_free_list;     // Free environment list
                                      // (linked by Env->env_link)
```

Listing 1: Three global variables pertaining environments defined in kern/env.c

For this three environment variables:

1. **envs** is a pointer pointing to an array of **Env** structures representing all the environments in the system.

2. **curenv** is used to keep track of the currently executing environment.

3. **env_free_list** is a array which keeps all of the inactive **Env** structures.

## 3.1 Environment State

The **Env** structure is defined in **inc/env.h** as follows:

```
struct Env {
        struct Trapframe env_tf;       // Saved registers
        struct Env *env_link;          // Next free Env
        envid_t env_id;                // Unique environment identifier
        envid_t env_parent_id;         // env_id of this env's parent
        enum EnvType env_type;         // Indicates special system environments
        unsigned env_status;           // Status of the environment
        uint32_t env_runs;             // Number of times environment has run

        // Address space
        pde_t *env_pgdir;              // Kernel virtual address of page dir
};
```

1. **env_tf**: This structure, defined in **inc/trap.h**, holds the saved register values for the environment while that environment is not running: i.e., when the kernel or a different environment is running. The kernel saves these when switching from user to kernel mode, so that the environment can later be resumed where it left off.

2. **env_link**: This is a link to the next **Env** on the **env_free_list**. **env_free_list** points to the first free environment on the list.

3. **env_id**: The kernel stores here a value that uniquely identifiers the environment currently using this Env structure (i.e., using this particular slot in the envs array). After a user environment terminates, the kernel may re-allocate the same Env structure to a different environment - but the new environment will have a different env_id from the old one even though the new environment is re-using the same slot in the envs array.

4. **env_parent_id**: The kernel stores here the env_id of the environment that created this environment. In this way the environments can form a family tree, which will be useful for making security decisions about which environments are allowed to do what to whom.

5. **env_type**: This is used to distinguish special environments. For most environments, it will be ENV_TYPE_USER. We'll introduce a few more types for special system service environments in later labs.

6. **env_status**: the status of this environment.

7. **env_pgdir**: holds the kernel virtual address of this environment's page directory.

## 3.2    Allocating the Environments Array

Modify **mem_init()** further to allocate a similar array of **Env** structures, called **envs**.

### 3.2.1    Exercise 1

*Modify **mem_init()** in **kern/pmap.c** to allocate and map the envs array. This array consists of exactly NENV instances of the **Env** structure allocated much like how you allocated the pages array. Also like the pages array, the memory backing envs should also be mapped user read-only at UENVS (defined in **inc/memlayout.h**) so user processes can read from this array.*

Firstly, let's check the hints in **mem_init()**

1. Hint 1:

```
        //////////////////////////////////////////////////////////////////////
        // Make 'envs' point to an array of size 'NENV' of 'struct Env'.
        // LAB 3: Your code here.
```

2. Hint 2:

```
        //////////////////////////////////////////////////////////////////////
        // Map the 'envs' array read-only by the user at linear address UENVS
        // (ie. perm = PTE_U | PTE_P).
        // Permissions:
        //    - the new image at UENVS  -- kernel R, user R
        //    - envs itself -- kernel RW, user NONE
        // LAB 3: Your code here.
```

For the above hints, we have our solutions for exercise 1.

```
        //////////////////////////////////////////////////////////////////////
        // Make 'envs' point to an array of size 'NENV' of 'struct Env'.
        // LAB 3: Your code here.
        envs = (struct Env*) boot_alloc(NENV*sizeof(struct Env));
        memset(envs, 0, NENV*sizeof(struct Env));
```

```
eee
```

## 3.3 Creating and Running Environments

### 3.3.1 Exercise 2

*In the file **env.c**, finish coding the following functions:*

1. ***env_init()***: *Initialize all of the Env structures in the **envs** array and add them to the **env_free_list**. Also calls **env_init_percpu**, which configures the segmentation hardware with separate segments for privilege level 0 (kernel) and privilege level 3 (user).*

2. ***env_setup_vm()***: *Allocate a page directory for a new environment and initialize the kernel portion of the new environment's address space.*

3. ***region_alloc()***: *Allocates and maps physical memory for an environment*

4. ***load_icode()***: *You will need to parse an ELF binary image, much like the boot loader already does, and load its contents into the user address space of a new environment.*

5. ***env_create()***: *Allocate an environment with **env_alloc** and call **load_icode** to load an ELF binary into it.*

6. ***env_run()***: *Start a given environment running in user mode.*

*As you write these functions, you might find the new cprintf verb %e useful – it prints a description corresponding to an error code. For example, **r = -E_NO_MEM; panic("env_alloc: %e", r);** will panic with the message "env_alloc: out of memory".*

### 3.3.2 The Implementation of env_init() Function

Let's firstly see the hint:

```
// Mark all environments in 'envs' as free, set their env_ids to 0,
// and insert them into the env_free_list.
// Make sure the environments are in the free list in the same order
// they are in the envs array (i.e., so that the first call to
// env_alloc() returns envs[0]).
//
void
env_init(void)
{
        // Set up envs array
        // LAB 3: Your code here.

        // Per-CPU part of the initialization
        env_init_percpu();
}
```

Next, see my implementation:

```
void
env_init(void)
{
        // Set up envs array
        // LAB 3: Your code here.
        env_free_list = NULL;
```

7

```
7         for (int i=NENV-1; i>=0; i--) {
8                 envs[i].env_id = 0;
9                 envs[i].env_status = ENV_FREE;
10                envs[i].env_link = env_free_list;
11                env_free_list = &envs[i];
12        }
13
14        // Per-CPU part of the initialization
15        env_init_percpu();
16 }
```

### 3.3.3 The Implementation of env_setup_vm() Function

This is for allocating a page directory for a new environment and only initialize the **kernal portion** of the new environment address space. Firstly, see the original code and its hints:

```
1  //
2  // Initialize the kernel virtual memory layout for environment e.
3  // Allocate a page directory, set e->env_pgdir accordingly,
4  // and initialize the kernel portion of the new environment's address space.
5  // Do NOT (yet) map anything into the user portion
6  // of the environment's virtual address space.
7  //
8  // Returns 0 on success, < 0 on error.  Errors include:
9  //      -E_NO_MEM if page directory or table could not be allocated.
10 //
11 static int
12 env_setup_vm(struct Env *e)
13 {
14        int i;
15        struct PageInfo *p = NULL;
16
17        // Allocate a page for the page directory
18        if (!(p = page_alloc(ALLOC_ZERO)))
19                return -E_NO_MEM;
20        // Now, set e->env_pgdir and initialize the page directory.
21        //
22        // Hint:
23        //    - The VA space of all envs is identical above UTOP
24        //      (except at UVPT, which we've set below).
25        //      See inc/memlayout.h for permissions and layout.
26        //      Can you use kern_pgdir as a template?  Hint: Yes.
27        //      (Make sure you got the permissions right in Lab 2.)
28        //    - The initial VA below UTOP is empty.
29        //    - You do not need to make any more calls to page_alloc.
30        //    - Note: In general, pp_ref is not maintained for
31        //      physical pages mapped only above UTOP, but env_pgdir
32        //      is an exception -- you need to increment env_pgdir's
33        //      pp_ref for env_free to work correctly.
34        //    - The functions in kern/pmap.h are handy.
35
36        // LAB 3: Your code here.
37
38        // UVPT maps the env's own page table read-only.
39        // Permissions: kernel R, user R
40        e->env_pgdir[PDX(UVPT)] = PADDR(e->env_pgdir) | PTE_P | PTE_U;
```

8

```
41
42          return 0;
43  }
```

```
1   static int
2   env_setup_vm(struct Env *e)
3   {
4           int i;
5           struct PageInfo *p = NULL;
6
7           // Allocate a page for the page directory
8           if (!(p = page_alloc(ALLOC_ZERO)))
9                   return -E_NO_MEM;
10          // LAB 3: Your code here.
11          // Step 1: set e->env_pgdir
12          e->env_pgdir = (pde_t *)page2kva(p);
13          p->pp_ref++;
14
15          // Step 2 and Step 3 is to initialize the page directory of this envir.
16          // Step 2: map the directory below UTOP
17          //         the initial VA below UTOP is empty.
18          for(i=0; i<PDX(UTOP); i++) {
19                  e->env_pgdir[i] = 0;
20          }
21          // Step 3: map the directory above UTOP
22          for(i=PDX(UTOP); i < NPDENTRIES; i++) {
23                  e->env_pgdir[i] = kern_pgdir[i];
24          }
25
26          // UVPT maps the env's own page table read-only.
27          // Permissions: kernel R, user R
28          e->env_pgdir[PDX(UVPT)] = PADDR(e->env_pgdir) | PTE_P | PTE_U;
29
30          return 0;
31  }
```

### 3.3.4   The Implementation of region_alloc() Function

**region_alloc()** function is used to allocate and map physical memory to environment.

```
1   //
2   // Allocate len bytes of physical memory for environment env,
3   // and map it at virtual address va in the environment's address space.
4   // Does not zero or otherwise initialize the mapped pages in any way.
5   // Pages should be writable by user and kernel.
6   // Panic if any allocation attempt fails.
7   //
8   static void
9   region_alloc(struct Env *e, void *va, size_t len)
10  {
11          // LAB 3: Your code here.
12          // (But only if you need it for load_icode.)
13          //
14          // Hint: It is easier to use region_alloc if the caller can pass
15          //    'va' and 'len' values that are not page-aligned.
16          //    You should round va down, and round (va + len) up.
17          //    (Watch out for corner-cases!)
```

```
18          struct PageInfo *p = NULL;
19          void *i;
20          int r;
21          void *start = (void *)ROUNDDOWN(va, PGSIZE);
22          void *end = (void *)ROUNDUP(va+len, PGSIZE);
23          for(i=start; i<end; i+=PGSIZE) {
24                  // page allocation
25                  p = page_alloc(0);
26                  if(p == NULL) {
27                          panic("region alloc: page alloc failed!");
28                  }
29                  // page insert
30                  // page insert
31                  r = page_insert(e->env_pgdir, p, i, PTE_W | PTE_U);
32                  if(r != 0)
33                          panic("region alloc: pgdir modification failed");
34          }
35  }
```

### 3.3.5   The Implementation of load_icode() Function

The initial hints:

```
1  //
2  // Set up the initial program binary, stack, and processor flags
3  // for a user process.
4  // This function is ONLY called during kernel initialization,
5  // before running the first user-mode environment.
6  //
7  // This function loads all loadable segments from the ELF binary image
8  // into the environment's user memory, starting at the appropriate
9  // virtual addresses indicated in the ELF program header.
10 // At the same time it clears to zero any portions of these segments
11 // that are marked in the program header as being mapped
12 // but not actually present in the ELF file - i.e., the program's bss section.
13 //
14 // All this is very similar to what our boot loader does, except the boot
15 // loader also needs to read the code from disk.  Take a look at
16 // boot/main.c to get ideas.
17 //
18 // Finally, this function maps one page for the program's initial stack.
19 //
20 // load_icode panics if it encounters problems.
21 //  - How might load_icode fail?  What might be wrong with the given input?
22 //
23 static void
24 load_icode(struct Env *e, uint8_t *binary)
25 {
26         // Hints:
27         //  Load each program segment into virtual memory
28         //  at the address specified in the ELF segment header.
29         //  You should only load segments with ph->p_type == ELF_PROG_LOAD.
30         //  Each segment's virtual address can be found in ph->p_va
31         //  and its size in memory can be found in ph->p_memsz.
32         //  The ph->p_filesz bytes from the ELF binary, starting at
33         //  'binary + ph->p_offset', should be copied to virtual address
34         //  ph->p_va.  Any remaining memory bytes should be cleared to zero.
```

```
35          //  (The ELF header should have ph->p_filesz <= ph->p_memsz.)
36          //  Use functions from the previous lab to allocate and map pages.
37          //
38          //  All page protection bits should be user read/write for now.
39          //  ELF segments are not necessarily page-aligned, but you can
40          //  assume for this function that no two segments will touch
41          //  the same virtual page.
42          //
43          //  You may find a function like region_alloc useful.
44          //
45          //  Loading the segments is much simpler if you can move data
46          //  directly into the virtual addresses stored in the ELF binary.
47          //  So which page directory should be in force during
48          //  this function?
49          //
50          //  You must also do something with the program's entry point,
51          //  to make sure that the environment starts executing there.
52          //  What?  (See env_run() and env_pop_tf() below.)
53
54          // LAB 3: Your code here.
55
56          // Now map one page for the program's initial stack
57          // at virtual address USTACKTOP - PGSIZE.
58
59          // LAB 3: Your code here.
60 }
```

**It's difficult for me!**

Two steps:

1. Parse an ELF binary image.

2. Load the ELP binary image content into the user address space of the new environment.

### 3.3.6 The Implementation of env_create() Function

```
1  // Allocates a new env with env_alloc, loads the named elf
2  // binary into it with load_icode, and sets its env_type.
3  // This function is ONLY called during kernel initialization,
4  // before running the first user-mode environment.
5  // The new env's parent ID is set to 0.
6  //
7  void
8  env_create(uint8_t *binary, enum EnvType type)
9  {
10         // LAB 3: Your code here.
11 }
```

There are mainly three steps:
1. Allocate a new env with **env_alloc()**.

2. Load the named elf binary into this env with **load_icode()**.

3. Set the env type.

It seems quite easy now:

```
1  void
2  env_create(uint8_t *binary, enum EnvType type)
3  {
4          // LAB 3: Your code here.
5          struct Env *env;
6          int rc;
7          // step 1: allocates a new env with env_alloc
8          rc = env_alloc(&e, 0);
9          if(rc != 0)
10                 panic("env_create: env_alloc failed");
11         // step 2: loads the named elf binary with load_icode
12         load_icode(env, binary);
13         // step 3: set env's env_type
14         env->env_type = type;
15
16 }
```

### 3.3.7   The Implementation of env_run() Function

Let's firstly revisit its hints:

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

Just follows the instructions:

```
1  void
2  env_run(struct Env *e)
3  {
4      // LAB 3: Your code here.
```

12

- **start** (kern/entry.S)
- **i386_init** (kern/init.c)
  - **cons_init**
  - **mem_init**
  - **env_init**
  - **trap_init** (still incomplete at this point)
  - **env_create**
  - **env_run**
    - **env_pop_tf**

Figure 1: Call graph

```
5         // Step 1: change the curenv's env_status from ING to ABLE.
6
7         if(curenv != NULL && curenv->env_status == ENV_RUNNING)
8                 curenv->env_status = ENV_RUNNABLE;
9
10        // Step 2: switch to the new env
11        curenv = e;
12        // Step 3: set status
13        curenv->env_status = ENV_RUNNING;
14        // Step 4: update env_runs counter
15        curenv->env_runs++;
16        // Step 5: use lcr3() to switch to its address space
17        lcr3(PADDR(curenv->env_pgdir));
18
19        // Step 6: use env_pop_tf() to restore the env's register
20        //         and drop into user mode in the environment.
21        env_pop_tf(&curenv->env_tf);
22
23
24        panic("env_run not yet implemented");
25 }
```

Let's check what's **lcr3()** function is for: It's defined in **inc/x86.h**

```
1 static inline void
2 lcr3(uint32_t val)
3 {
4         asm volatile("movl %0,%%cr3" : : "r" (val));
5 }
```

## 3.4   The Call Graph

## 3.5   Basics of Protected Control Transfer

**Exceptions** and **interrupts** are both "protected control transfers," which cause the processor to switch from *user* to *kernel* mode (CPL=0) without giving the user-mode

13

code any opportunity to interfere with the functioning of the kernel or other environments.

There is difference between **exceptions** and **interrupts** in Intel's terminology:

- **An interrupt** is a protected control transfer that is caused by an asynchronous event usually external to the processor, such as notification of external device I/O activity.

- **An exception**, in contrast, is a protected control transfer caused synchronously by the currently running code, for example due to a divide by zero or an invalid memory access

There are two mechanism working together to provide the control transfer protection:

1. **The Interrupt Descriptor Table (IDT):**

2. **The Task State Segment:** The processor ensures that interrupts and exceptions can only cause the kernel to be entered at a few specific, well-defined entry-points determined by the kernel itself, and not by the code running when the interrupt or exception is taken. **The Task State Segment:** The processor needs a place to save the old processor state before the interrupt or exception occurred, such as the original values of EIP and CS before the processor invoked the exception handler, so that the exception handler can later restore that old state and resume the interrupted code from where it left off. But this save area for the old processor state must in turn be protected from unprivileged user-mode code; otherwise buggy or malicious user code could compromise the kernel.

## 3.6   Types of Exceptions and Interrupts

1. **Synchronous exceptions (the processor exception)**: All of the synchronous exceptions that the x86 processor can generate internally use interrupt vectors between 0 and 31, and therefore map to IDT entries 0-31.

2. **Software interrupts and Asynchronous hardware interrupts**: interrupt vectors greater than 31 are only used by software interrupts, which can be generated by the int instruction or asynchronous hardware interrupts, caused by external devices when they need attention.

## 3.7   Setting Up the IDT

Let's see the details:

1. The **trap_init()** function initialize the IDT with the address of exception or interrupt handlers.

2. These handlers are implemented in **trapentry.S**.

3. Each of these handlers should build a **struct Trapframe** (see in **inc/trap.h**) on the stack and call **trap()** with a pointer to the Trapframe.

4. **trap()** then handles the exception/interrupt or dispatches to a specific handler function.

## 3.8   IDT of x86

# Interrupt/trap gate

31     16 15 14 13 12    8 7   5 4     0

| Offset 31..16 | P | D P L | 0 D 1 1 | **Type** | 0 0 0 | |

31     16 15     0

| Segment selector | Offset 15..0 |

Type     0=Interrupt gate        P    Present
              1=Trap gate           DPL Descriptor privilege level
Selector Destination CS              (CPL required to invoke gate)
Offset    Destination IP or EIP    D    Size of gate (0=16-bits, 1=32-bits)

## Exception stack (with privilege change)

| | ← SS:ESP from TSS |
| SS | |
| ESP | |
| EFLAGS | |
| CS | |
| EIP | |
| Error code | ← ESP at handler entry |

## Exception stack (without privilege change)

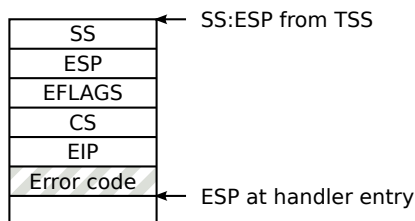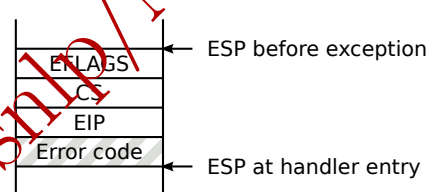| | ← ESP before exception |
| EFLAGS | |
| CS | |
| EIP | |
| Error code | ← ESP at handler entry |

| Vector | Description | Type | Error code | Exception types | | |
|---|---|---|---|---|---|---|
| 0 | Divide error | Fault | No | Fault | Faulting instruction not executed CS:EIP is the faulting instruction | |
| 1 | Reserved | | | | | |
| 2 | Non-maskable interrupt | Interrupt | No | Trap | Trapping instruction executed CS:EIP is the next instruction | |
| 3 | Breakpoint | Trap | No | | | |
| 4 | Overflow | Trap | No | Abort | Location is imprecise; cannot safely resume execution | |
| 5 | BOUND range exceeded | Fault | No | | | |
| 6 | Invalid/undefined opcode | Fault | No | | | |
| 7 | No math coprocessor | Fault | No | | | |
| 8 | Double fault | Abort | Zero | | | |
| 9 | Reserved | | | | | |
| 10 | Invalid TSS | Fault | Yes | | | |
| 11 | Segment not present | Fault | Yes | TI | 0=GDT, 1=LDT | |
| 12 | Stack-segment fault | Fault | Yes | IDT | 0=GDT/LDT, 1=IDT | |
| 13 | General protection | Fault | Yes | EXT | External event | |
| 14 | Page fault | Fault | Yes | | | |
| 15 | Reserved | | No | | | |
| 16 | x87 FPU error | Fault | No | | | |
| 17 | Alignment check | Fault | Zero | | | |
| 18 | Machine check | Abort | No | P | 0=Non-present page 1=Protection-violation | |
| 19 | SIMD FP exception | Fault | No | W/R | Cause (0=Read, 1=Write) | |
| 20-31 | Reserved | | | U/S | Mode (0=Supervisor, 1=User) | |
| 32-255 | User defined interrupts | Interrupt | No | | | |

31     3 2 1 0

| Segment selector | TI | IDT | EXT |

31     3 2 1 0

| | U/S | W/R | P |

```
           IDT                    trapentry.S          trap.c
       +---------------+
       |    &handler1  |---------> handler1:           trap (struct Trapframe *tf)
       |               |              // do stuff       {
       |               |              call trap           // handle the exception/interrupt
       |               |              // ...            }
       +---------------+
       |    &handler2  |--------> handler2:
       |               |             // do stuff
       |               |             call trap
       |               |             // ...
       +---------------+
            .
            .
            .
       +---------------+
       |    &handlerX  |--------> handlerX:
       |               |             // do stuff
       |               |             call trap
       |               |             // ...
       +---------------+
```
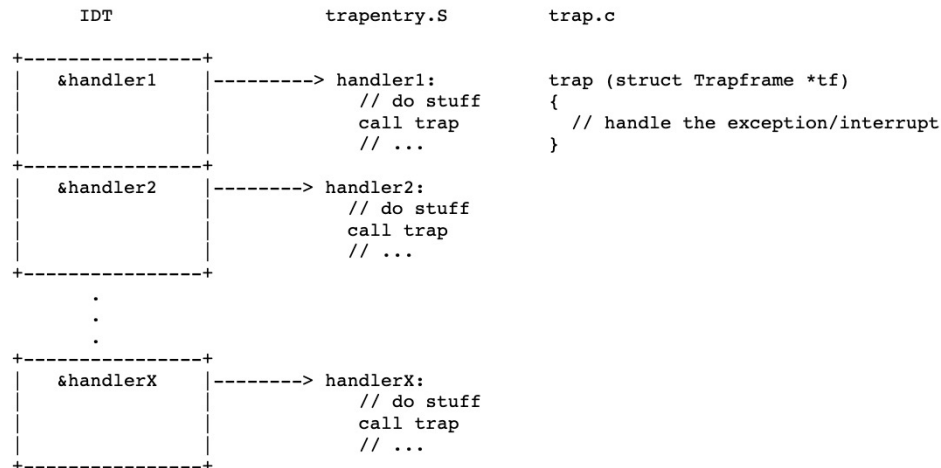
Figure 2: Setting up the IDT

**Exercise 4.** Edit `trapentry.s` and `trap.c` and implement the features described above. The macros `TRAPHANDLER` and `TRAPHANDLER_NOEC` in `trapentry.s` should help you, as well as the T_* defines in `inc/trap.h`. You will need to add an entry point in `trapentry.s` (using those macros) for each trap defined in `inc/trap.h`, and you'll have to provide `_alltraps` which the `TRAPHANDLER` macros refer to. You will also need to modify `trap_init()` to initialize the `idt` to point to each of these entry points defined in `trapentry.s`; the `SETGATE` macro will be helpful here.

Your `_alltraps` should:

1. push values to make the stack look like a struct Trapframe
2. load `GD_KD` into `%ds` and `%es`
3. `pushl %esp` to pass a pointer to the Trapframe as an argument to trap()
4. `call trap` (can `trap` ever return?)

Consider using the `pushal` instruction; it fits nicely with the layout of the `struct Trapframe`.

Test your trap handling code using some of the test programs in the `user` directory that cause exceptions before making any system calls, such as `user/divzero`. You should be able to get **make grade** to succeed on the `divzero`, `softint`, and `badsegment` tests at this point.

Figure 3: Exercise 4

## 3.9   Exercise 4

### 3.9.1   The Hints of kern/trapentry.S

```
1  /* See COPYRIGHT for copyright information. */
2
3  #include <inc/mmu.h>
4  #include <inc/memlayout.h>
5  #include <inc/trap.h>
6
7
8
9  ###################################################################
10 # exceptions/interrupts
11 ###################################################################
12
```

16

```
13  /* TRAPHANDLER defines a globally-visible function for handling a trap.
14   * It pushes a trap number onto the stack, then jumps to _alltraps.
15   * Use TRAPHANDLER for traps where the CPU automatically pushes an error code.
16   *
17   * You shouldn't call a TRAPHANDLER function from C, but you may
18   * need to _declare_ one in C (for instance, to get a function pointer
19   * during IDT setup).  You can declare the function with
20   *   void NAME();
21   * where NAME is the argument passed to TRAPHANDLER.
22   */
23  #define TRAPHANDLER(name, num)            \
24    .globl name;    /* define global symbol for 'name' */ \
25    .type name, @function;  /* symbol type is function */   \
26    .align 2;   /* align function definition */   \
27    name:     /* function starts here */    \
28    pushl $(num);            \
29    jmp _alltraps
30
31  /* Use TRAPHANDLER_NOEC for traps where the CPU doesn't push an error code.
32   * It pushes a 0 in place of the error code, so the trap frame has the same
33   * format in either case.
34   */
35  #define TRAPHANDLER_NOEC(name, num)          \
36    .globl name;            \
37    .type name, @function;          \
38    .align 2;          \
39    name:            \
40    pushl $0;          \
41    pushl $(num);          \
42    jmp _alltraps
43
44  .text
45
46  /*
47   * Lab 3: Your code here for generating entry points for the different traps.
48   */
49
50
51
52  /*
53   * Lab 3: Your code here for _alltraps
54   */
```

We should notice that macros **TRAPHANDLER** and **TRAPHANDLER_NOEC** (no error code needs to be pushed) are used for traps. That's enough.

Next, let's see the content in file inc/trap.h

```
1  #ifndef JOS_INC_TRAP_H
2  #define JOS_INC_TRAP_H
3
4  // Trap numbers
5  // These are processor defined:
6  #define T_DIVIDE    0    // divide error
7  #define T_DEBUG     1    // debug exception
8  #define T_NMI       2    // non-maskable interrupt
9  #define T_BRKPT     3    // breakpoint
10 #define T_OFLOW     4    // overflow
11 #define T_BOUND     5    // bounds check
```

```c
#define T_ILLOP      6    // illegal opcode
#define T_DEVICE     7    // device not available
#define T_DBLFLT     8    // double fault
/* #define T_COPROC  9 */ // reserved (not generated by recent processors)
#define T_TSS        10   // invalid task switch segment
#define T_SEGNP      11   // segment not present
#define T_STACK      12   // stack exception
#define T_GPFLT      13   // general protection fault
#define T_PGFLT      14   // page fault
/* #define T_RES     15 */ // reserved
#define T_FPERR      16   // floating point error
#define T_ALIGN      17   // aligment check
#define T_MCHK       18   // machine check
#define T_SIMDERR    19   // SIMD floating point error

// These are arbitrarily chosen, but with care not to overlap
// processor defined exceptions or interrupt vectors.
#define T_SYSCALL    48   // system call
#define T_DEFAULT    500  // catchall

#define IRQ_OFFSET  32  // IRQ 0 corresponds to int IRQ_OFFSET

// Hardware IRQ numbers. We receive these as (IRQ_OFFSET+IRQ_WHATEVER)
#define IRQ_TIMER       0
#define IRQ_KBD         1
#define IRQ_SERIAL      4
#define IRQ_SPURIOUS    7
#define IRQ_IDE         14
#define IRQ_ERROR       19

#ifndef __ASSEMBLER__

#include <inc/types.h>

struct PushRegs {
  /* registers as pushed by pusha */
  uint32_t reg_edi;
  uint32_t reg_esi;
  uint32_t reg_ebp;
  uint32_t reg_oesp;      /* Useless */
  uint32_t reg_ebx;
  uint32_t reg_edx;
  uint32_t reg_ecx;
  uint32_t reg_eax;
} __attribute__((packed));

struct Trapframe {
  struct PushRegs tf_regs;
  uint16_t tf_es;
  uint16_t tf_padding1;
  uint16_t tf_ds;
  uint16_t tf_padding2;
  uint32_t tf_trapno;
  /* below here defined by x86 hardware */
  uint32_t tf_err;
  uintptr_t tf_eip;
  uint16_t tf_cs;
  uint16_t tf_padding3;
```

```
70   uint32_t tf_eflags;
71   /* below here only when crossing rings, such as from user to kernel */
72   uintptr_t tf_esp;
73   uint16_t tf_ss;
74   uint16_t tf_padding4;
75 } __attribute__((packed));
76
77
78 #endif /* !__ASSEMBLER__ */
79
80 #endif /* !JOS_INC_TRAP_H */
```

### 3.9.2   The Hints of kern/trap.c

Actually, we should also see the code in [kern/trap.c](kern/trap.c).

```
1  #include <inc/mmu.h>
2  #include <inc/x86.h>
3  #include <inc/assert.h>
4
5  #include <kern/pmap.h>
6  #include <kern/trap.h>
7  #include <kern/console.h>
8  #include <kern/monitor.h>
9  #include <kern/env.h>
10 #include <kern/syscall.h>
11
12 static struct Taskstate ts;
13
14 /* For debugging, so print_trapframe can distinguish between printing
15  * a saved trapframe and printing the current trapframe and print some
16  * additional information in the latter case.
17  */
18 static struct Trapframe *last_tf;
19
20 /* Interrupt descriptor table.  (Must be built at run time because
21  * shifted function addresses can't be represented in relocation records.)
22  */
23 struct Gatedesc idt[256] = { { 0 } };
24 struct Pseudodesc idt_pd = {
25   sizeof(idt) - 1, (uint32_t) idt
26 };
27
28
29 static const char *trapname(int trapno)
30 {
31   static const char * const excnames[] = {
32     "Divide error",
33     "Debug",
34     "Non-Maskable Interrupt",
35     "Breakpoint",
36     "Overflow",
37     "BOUND Range Exceeded",
38     "Invalid Opcode",
39     "Device Not Available",
40     "Double Fault",
41     "Coprocessor Segment Overrun",
42     "Invalid TSS",
43     "Segment Not Present",
```

```
44        "Stack Fault",
45        "General Protection",
46        "Page Fault",
47        "(unknown trap)",
48        "x87 FPU Floating-Point Error",
49        "Alignment Check",
50        "Machine-Check",
51        "SIMD Floating-Point Exception"
52      };
53
54      if (trapno < ARRAY_SIZE(excnames))
55        return excnames[trapno];
56      if (trapno == T_SYSCALL)
57        return "System call";
58      return "(unknown trap)";
59    }
60
61
62    void
63    trap_init(void)
64    {
65      extern struct Segdesc gdt[];
66
67      // LAB 3: Your code here.
68
69      // Per-CPU setup
70      trap_init_percpu();
71    }
72
73    // Initialize and load the per-CPU TSS and IDT
74    void
75    trap_init_percpu(void)
76    {
77      // Setup a TSS so that we get the right stack
78      // when we trap to the kernel.
79      ts.ts_esp0 = KSTACKTOP;
80      ts.ts_ss0 = GD_KD;
81      ts.ts_iomb = sizeof(struct Taskstate);
82
83      // Initialize the TSS slot of the gdt.
84      gdt[GD_TSS0 >> 3] = SEG16(STS_T32A, (uint32_t) (&ts),
85              sizeof(struct Taskstate) - 1, 0);
86      gdt[GD_TSS0 >> 3].sd_s = 0;
87
88      // Load the TSS selector (like other segment selectors, the
89      // bottom three bits are special; we leave them 0)
90      ltr(GD_TSS0);
91
92      // Load the IDT
93      lidt(&idt_pd);
94    }
95
96    void
97    print_trapframe(struct Trapframe *tf)
98    {
99      cprintf("TRAP frame at %p\n", tf);
100     print_regs(&tf->tf_regs);
101     cprintf("  es   0x----%04x\n", tf->tf_es);
```

```
102    cprintf("  ds   0x----%04x\n", tf->tf_ds);
103    cprintf("  trap 0x%08x %s\n", tf->tf_trapno, trapname(tf->tf_trapno));
104    // If this trap was a page fault that just happened
105    // (so %cr2 is meaningful), print the faulting linear address.
106    if (tf == last_tf && tf->tf_trapno == T_PGFLT)
107      cprintf("  cr2  0x%08x\n", rcr2());
108    cprintf("  err  0x%08x", tf->tf_err);
109    // For page faults, print decoded fault error code:
110    // U/K=fault occurred in user/kernel mode
111    // W/R=a write/read caused the fault
112    // PR=a protection violation caused the fault (NP=page not present).
113    if (tf->tf_trapno == T_PGFLT)
114      cprintf(" [%s, %s, %s]\n",
115        tf->tf_err & 4 ? "user" : "kernel",
116        tf->tf_err & 2 ? "write" : "read",
117        tf->tf_err & 1 ? "protection" : "not-present");
118    else
119      cprintf("\n");
120    cprintf("  eip  0x%08x\n", tf->tf_eip);
121    cprintf("  cs   0x----%04x\n", tf->tf_cs);
122    cprintf("  flag 0x%08x\n", tf->tf_eflags);
123    if ((tf->tf_cs & 3) != 0) {
124      cprintf("  esp  0x%08x\n", tf->tf_esp);
125      cprintf("  ss   0x----%04x\n", tf->tf_ss);
126    }
127  }
128
129  void
130  print_regs(struct PushRegs *regs)
131  {
132    cprintf("  edi  0x%08x\n", regs->reg_edi);
133    cprintf("  esi  0x%08x\n", regs->reg_esi);
134    cprintf("  ebp  0x%08x\n", regs->reg_ebp);
135    cprintf("  oesp 0x%08x\n", regs->reg_oesp);
136    cprintf("  ebx  0x%08x\n", regs->reg_ebx);
137    cprintf("  edx  0x%08x\n", regs->reg_edx);
138    cprintf("  ecx  0x%08x\n", regs->reg_ecx);
139    cprintf("  eax  0x%08x\n", regs->reg_eax);
140  }
141
142  static void
143  trap_dispatch(struct Trapframe *tf)
144  {
145    // Handle processor exceptions.
146    // LAB 3: Your code here.
147
148    // Unexpected trap: The user process or the kernel has a bug.
149    print_trapframe(tf);
150    if (tf->tf_cs == GD_KT)
151      panic("unhandled trap in kernel");
152    else {
153      env_destroy(curenv);
154      return;
155    }
156  }
157
158  void
159  trap(struct Trapframe *tf)
```

21

```
160 {
161   // The environment may have set DF and some versions
162   // of GCC rely on DF being clear
163   asm volatile("cld" ::: "cc");
164
165   // Check that interrupts are disabled.  If this assertion
166   // fails, DO NOT be tempted to fix it by inserting a "cli" in
167   // the interrupt path.
168   assert(!(read_eflags() & FL_IF));
169
170   cprintf("Incoming TRAP frame at %p\n", tf);
171
172   if ((tf->tf_cs & 3) == 3) {
173     // Trapped from user mode.
174     assert(curenv);
175
176     // Copy trap frame (which is currently on the stack)
177     // into 'curenv->env_tf', so that running the environment
178     // will restart at the trap point.
179     curenv->env_tf = *tf;
180     // The trapframe on the stack should be ignored from here on.
181     tf = &curenv->env_tf;
182   }
183
184   // Record that tf is the last real trapframe so
185   // print_trapframe can print some additional information.
186   last_tf = tf;
187
188   // Dispatch based on what type of trap occurred
189   trap_dispatch(tf);
190
191   // Return to the current environment, which should be running.
192   assert(curenv && curenv->env_status == ENV_RUNNING);
193   env_run(curenv);
194 }
195
196
197 void
198 page_fault_handler(struct Trapframe *tf)
199 {
200   uint32_t fault_va;
201
202   // Read processor's CR2 register to find the faulting address
203   fault_va = rcr2();
204
205   // Handle kernel-mode page faults.
206
207   // LAB 3: Your code here.
208
209   // We've already handled kernel-mode exceptions, so if we get here,
210   // the page fault happened in user mode.
211
212   // Destroy the environment that caused the fault.
213   cprintf("[%08x] user fault va %08x ip %08x\n",
214     curenv->env_id, fault_va, tf->tf_eip);
215   print_trapframe(tf);
216   env_destroy(curenv);
217 }
```

**kern/trap.c** has these functions:

1. **static const char \*trapname(int trapno)** : get trap name.

2. **void trap_init(void)**: trap initialization.

3. **void trap_init_percpu(void)**: Initialize and load the per-CPU TSS and IDT, called from **trap_init**.

4. **void print_trapframe(struct Trapframe \*tf)**: as what the name suggests.

5. **void print_regs(struct PushRegs \*regs)**: print registers' value.

6. **static void trap_dispatch(struct Trapframe \*tf)**: dispatch the traps according to the trapframe.

7. **void trap(struct Trapframe \*tf)**: the main trap function.

8. **void page_fault_handler(struct Trapframe \*tf)**

### 3.9.3  The Implementation of kern/trapentry.S

```
1  /*
2   * Lab 3: Your code here for generating entry points for the different traps.
3   */
4  TRAPHANDLER_NOEC(DIVIDE, T_DIVIDE)
5  TRAPHANDLER_NOEC(DEBUG, T_DEBUG)
6  TRAPHANDLER_NOEC(NMI, T_NMI)
7  TRAPHANDLER_NOEC(BRKPT, T_BRKPT)
8  TRAPHANDLER_NOEC(OFLOW, T_OFLOW)
9  TRAPHANDLER_NOEC(BOUND, T_BOUND)
10 TRAPHANDLER_NOEC(ILLOP, T_ILLOP)
11 TRAPHANDLER_NOEC(DEVICE, T_DEVICE)
12 TRAPHANDLER(DBLFLT, T_DBLFLT)
13 TRAPHANDLER(TSS, T_TSS)
14 TRAPHANDLER(SEGNP, T_SEGNP)
15 TRAPHANDLER(STACK, T_STACK)
16 TRAPHANDLER(GPFLT, T_GPFLT)
17 TRAPHANDLER(PGFLT, T_PGFLT)
18 TRAPHANDLER_NOEC(GPFLT, T_GPFLT)
19 TRAPHANDLER(ALIGN, T_ALIGN)
20 TRAPHANDLER_NOEC(MCHK, T_MCHK)
21 TRAPHANDLER_NOEC(SIMDERR, T_SIMDERR)
22 /*
23  SYSCALL: 48
24 */
25 TRAPHANDLER_NOEC(SYSCALL, T_SYSCALL)
26
27
28 /*
29  * Lab 3: Your code here for _alltraps
30  */
31 _alltraps:
32   pushl %ds
33   pushl %es
34   pushall
35
```

```
36    movl $GD_KD, %eax
37    movw %ax, %ds
38    movw %ax, %es
39
40    push %esp
41    call trap
```

### 3.9.4   The Implementation trap_init() in kern/trap.c

Let's firstly see the macros **SETGATE** in inc/mmu.h

```
1  // Set up a normal interrupt/trap gate descriptor.
2  // - istrap: 1 for a trap (= exception) gate, 0 for an interrupt gate.
3  //    see section 9.6.1.3 of the i386 reference: "The difference between
4  //    an interrupt gate and a trap gate is in the effect on IF (the
5  //    interrupt-enable flag). An interrupt that vectors through an
6  //    interrupt gate resets IF, thereby preventing other interrupts from
7  //    interfering with the current interrupt handler. A subsequent IRET
8  //    instruction restores IF to the value in the EFLAGS image on the
9  //    stack. An interrupt through a trap gate does not change IF."
10 // - sel: Code segment selector for interrupt/trap handler
11 // - off: Offset in code segment for interrupt/trap handler
12 // - dpl: Descriptor Privilege Level -
13 //        the privilege level required for software to invoke
14 //        this interrupt/trap gate explicitly using an int instruction.
15 #define SETGATE(gate, istrap, sel, off, dpl)                  \
16 {                                                             \
17         (gate).gd_off_15_0 = (uint32_t) (off) & 0xffff;       \
18         (gate).gd_sel = (sel);                                \
19         (gate).gd_args = 0;                                   \
20         (gate).gd_rsv1 = 0;                                   \
21         (gate).gd_type = (istrap) ? STS_TG32 : STS_IG32;      \
22         (gate).gd_s = 0;                                      \
23         (gate).gd_dpl = (dpl);                                \
24         (gate).gd_p = 1;                                      \
25         (gate).gd_off_31_16 = (uint32_t) (off) >> 16;         \
26 }
```

## 3.10   Questions Following Exercise 4

1. *What is the purpose of having an individual handler function for each exception/interrupt? (i.e., if all exceptions/interrupts were delivered to the same handler, what feature that exists in the current implementation could not be provided?)*

2. *Did you have to do anything to make the **user/softint** program behave correctly? The grade script expects it to produce a general protection fault (trap 13), but **softint**'s code says **int $14**. Why should this produce interrupt vector 13? What happens if the kernel actually allows **softint's int $14** instruction to invoke the kernel's page fault handler (which is interrupt vector 14)?*

# 4 Part B: Page Faults, Breakpoints Exceptions, and System Calls

## 4.1 Handling Page Faults

The page fault exception: interrupt vector 14(T_PGFLT). When the processor takes a page fault, it stores the linear (virtual address) that caused this fault and store it in a special register: **CR2**.

### 4.1.1 Exercise 5

*Modify* **trap_dispatch()** *to dispatch page fault exceptions to* **page_fault_handler()**. *You should now be able to get* **make grade** *to succeed on the* **faultread, faultread-kernel, faultwrite, and faultwritekernel** *tests. If any of them don't work, figure out why and fix them. Remember that you can boot JOS into a particular user program using* **make run-x** *or* **make run-x-nox**. *For instance,* **make run-hello-nox** *runs the hello user program.*

The hints of original **trap_dispatch()**

```
static void
trap_dispatch(struct Trapframe *tf)
{
        // Handle processor exceptions.
        // LAB 3: Your code here.

        // Unexpected trap: The user process or the kernel has a bug.
        print_trapframe(tf);
        if (tf->tf_cs == GD_KT)
                panic("unhandled trap in kernel");
        else {
                env_destroy(curenv);
                return;
        }
}
```

The implementation is quite simple:

```
static void
trap_dispatch(struct Trapframe *tf)
{
        // Handle processor exceptions.
        // LAB 3: Your code here.
        switch(tf-tf_trapno) {
                case T_PGFLT:
                        page_fault_handler(tf);
                default:
                        break;
        }

        // Unexpected trap: The user process or the kernel has a bug.
        print_trapframe(tf);
        if (tf->tf_cs == GD_KT)
                panic("unhandled trap in kernel");
        else {
                env_destroy(curenv);
                return;
```

```
20          }
21  }
```

## 4.2   The Breakpoint Exception

### 4.2.1   Exercise 6

*Modify trap_dispatch() to make breakpoint exceptions invoke the kernel monitor. You should now be able to get make grade to succeed on the breakpoint test.*

The implementation is quite simple:

```
1          case T_BRKPT:
2                  monitor(tf);
3                  return;
```

But it seems do not work!!! Why, recall that we use

```
1  SETGATE(idt[T_BRKPT], 0, GD_KT, t_brkpt, 0);
```

However, the truely breakpoint use **int \$3**. The privilege level doesn't match.

So, we change the above command in **trap_init()** to

```
1  SETGATE(idt[T_BRKPT], 0, GD_KT, t_brkpt, 3);
```

## 4.3   System Calls: Exercise 7

User processes ask the kernel to do things for them by invoking system calls. When the user process invokes a system call, the processor enters kernel mode, the processor and the kernel cooperate to save the user process's state, the kernel executes appropriate code in order to carry out the system call, and then resumes the user process.

### 4.3.1   Exercise 7

**Exercise 7.** Add a handler in the kernel for interrupt vector `T_SYSCALL`. You will have to edit `kern/trapentry.S` and `kern/trap.c`'s `trap_init()`. You also need to change `trap_dispatch()` to handle the system call interrupt by calling `syscall()` (defined in `kern/syscall.c`) with the appropriate arguments, and then arranging for the return value to be passed back to the user process in `%eax`. Finally, you need to implement `syscall()` in `kern/syscall.c`. Make sure `syscall()` returns `-E_INVAL` if the system call number is invalid. You should read and understand `lib/syscall.c` (especially the inline assembly routine) in order to confirm your understanding of the system call interface. Handle all the system calls listed in `inc/syscall.h` by invoking the corresponding kernel function for each call.

Run the `user/hello` program under your kernel (`make run-hello`). It should print "hello, world" on the console and then cause a page fault in user mode. If this does not happen, it probably means your system call handler isn't quite right. You should also now be able to get `make grade` to succeed on the `testbss` test.

Figure 4: Exercise 7

The whole procedure is as follows:

1. **trap dispatch()** handle the **T SYSCALL**.

2. **trap dispatch()** calls **syscall()** in **kern/syscall.c**

See the hints about **syscall()** in **kern/syscall.c**:

```
// Dispatches to the correct kernel function, passing the arguments.
int32_t
syscall(uint32_t syscallno, uint32_t a1, uint32_t a2, uint32_t a3, uint32_t a4,
    ↪ uint32_t a5)
{
        // Call the function corresponding to the 'syscallno' parameter.
        // Return any appropriate return value.
        // LAB 3: Your code here.

        panic("syscall not implemented");

        switch (syscallno) {
        default:
                return -E_INVAL;
        }
}
```

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

```
// Dispatches to the correct kernel function, passing the arguments.
int32_t
syscall(uint32_t syscallno, uint32_t a1, uint32_t a2, uint32_t a3, uint32_t a4,
    ↪ uint32_t a5)
{
        // Call the function corresponding to the 'syscallno' parameter.
        // Return any appropriate return value.
        // LAB 3: Your code here.


        //panic("syscall not implemented");

        switch (syscallno) {
        case SYS_cputs:
                sys_cputs((char *)a1, a2);
                return 0;
        case SYS_cgetc:
                return sys_cgetc();
        case SYS_getenvid:
                return sys_getenvid();
        case SYS_env_destroy:
                return sys_env_destroy(a1);
        default:
                return -E_INVAL;
        }
}
```

**syscall()** calls **sys cputs**, which also requires implementation.

### 4.3.3 The Implementation of sys_cputs() in kern/syscall.c

```
// Print a string to the system console.
// The string is exactly 'len' characters long.
```

27

```
 3  // Destroys the environment on memory errors.
 4  static void
 5  sys_cputs(const char *s, size_t len)
 6  {
 7          // Check that the user has permission to read memory [s, s+len).
 8          // Destroy the environment if not.
 9
10          // LAB 3: Your code here.
11          if (curenv->env_tf.tf_cs & 3)
12                  user_mem_assert(curenv, s, len, 0);
13
14          // Print the string supplied by the user.
15          cprintf("%.*s", len, s);
16  }
```

### 4.3.4  DO NOT FORGET TRAP_INIT and TRAP_DISPATCH

In **trap_init()**:

```
1          void SYSCALL();
2          SETGATE(idt[48], 0, GD_KT, SYSCALL, 3);
```

In **trap_dispatch()**,

```
 1  static void
 2  trap_dispatch(struct Trapframe *tf)
 3  {
 4          // Handle processor exceptions.
 5          // LAB 3: Your code here.
 6          switch(tf->tf_trapno) {
 7                  case T_PGFLT:
 8                          page_fault_handler(tf);
 9                          return;
10                  case T_BRKPT:
11                          monitor(tf);
12                          return;
13                  case T_SYSCALL:
14                          tf->tf_regs.reg_eax = syscall(tf->tf_regs.reg_eax,
15                                  tf->tf_regs.reg_edx,
16                                  tf->tf_regs.reg_ecx, tf->tf_regs.reg_ebx,
17                                  tf->tf_regs.reg_edi, tf->tf_regs.reg_esi);
18                          return;
19                  default:
20                          break;
21          }
22
23          // Unexpected trap: The user process or the kernel has a bug.
24          print_trapframe(tf);
25          if (tf->tf_cs == GD_KT)
26                  panic("unhandled trap in kernel");
27          else {
28                  env_destroy(curenv);
29                  return;
30          }
31  }
```

NOW: 60/80 score.

## 4.4 User-mode startup

A user program starts from **lib/entry.S**. After some setup, this code calls **libmain()** in **lib/libmain.c**.

### 4.4.1 Exercise 8

A user program starts running at the top of `lib/entry.s`. After some setup, this code calls `libmain()`, in `lib/libmain.c`. You should modify `libmain()` to initialize the global pointer `thisenv` to point at this environment's `struct Env` in the `envs[]` array. (Note that `lib/entry.s` has already defined `envs` to point at the UENVS mapping you set up in Part A.) Hint: look in `inc/env.h` and use `sys_getenvid`.

`libmain()` then calls `umain`, which, in the case of the hello program, is in `user/hello.c`. Note that after printing "`hello, world`", it tries to access `thisenv->env_id`. This is why it faulted earlier. Now that you've initialized `thisenv` properly, it should not fault. If it still faults, you probably haven't mapped the UENVS area user–readable (back in Part A in `pmap.c`; this is the first time we've actually used the UENVS area).

> **Exercise 8.** Add the required code to the user library, then boot your kernel. You should see `user/hello` print "`hello, world`" and then print "`i am environment 00001000`". `user/hello` then attempts to "exit" by calling `sys_env_destroy()` (see `lib/libmain.c` and `lib/exit.c`). Since the kernel currently only supports one user environment, it should report that it has destroyed the only environment and then drop into the kernel monitor. You should be able to get **make grade** to succeed on the `hello` test.

Figure 5: Exercise 8

### 4.4.2 The Content in lib/entry.S

```
1   #include <inc/mmu.h>
2   #include <inc/memlayout.h>
3
4   .data
5   // Define the global symbols 'envs', 'pages', 'uvpt', and 'uvpd'
6   // so that they can be used in C as if they were ordinary global arrays.
7     .globl envs
8     .set envs, UENVS
9     .globl pages
10    .set pages, UPAGES
11    .globl uvpt
12    .set uvpt, UVPT
13    .globl uvpd
14    .set uvpd, (UVPT+(UVPT>>12)*4)
15
16
17  // Entrypoint - this is where the kernel (or our parent environment)
18  // starts us running when we are initially loaded into a new environment.
19  .text
20  .globl _start
21  _start:
22    // See if we were started with arguments on the stack
23    cmpl $USTACKTOP, %esp
24    jne args_exist
25
26    // If not, push dummy argc/argv arguments.
```

```
27    // This happens when we are loaded by the kernel,
28    // because the kernel does not know about passing arguments.
29    pushl $0
30    pushl $0
31
32 args_exist:
33    call libmain
34 1:  jmp 1b
```

### 4.4.3   libmain() Function in lib/libmain.c

```
1  // Called from entry.S to get us going.
2  // entry.S already took care of defining envs, pages, uvpd, and uvpt.
3
4  #include <inc/lib.h>
5
6  extern void umain(int argc, char **argv);
7
8  const volatile struct Env *thisenv;
9  const char *binaryname = "<unknown>";
10
11 void
12 libmain(int argc, char **argv)
13 {
14         // set thisenv to point at our Env structure in envs[].
15         // LAB 3: Your code here.
16         thisenv = &envs[ENVX(sys_getenvid())];
17
18         // save the name of the program so that panic() can use it
19         if (argc > 0)
20                 binaryname = argv[0];
21
22         // call user main routine
23         umain(argc, argv);
24
25         // exit gracefully
26         exit();
27 }
```

Now: 65/80.

## 4.5   Page faults and memory protection

### 4.5.1   Crucial Memory Protection Challenges Brought By System Calls

Most system call interfaces let user programs pass pointers to the kernel. These pointers
point at user buffers to be read or written. The kernel then dereferences these pointers
while carrying out the system call. There are two problems with this:

1. A page fault in the kernel is potentially a lot more serious than a page fault in a
   user program. If the kernel page-faults while manipulating its own data structures,
   that's a kernel bug, and the fault handler should panic the kernel (and hence the
   whole system). **But when the kernel is dereferencing pointers given to it
   by the user program, it needs a way to remember that any page faults
   these dereferences cause are actually on behalf of the user program.**

2. The kernel typically has more memory permissions than the user program. The
   user program might pass a pointer to a system call that points to memory that

the kernel can read or write but that the program cannot. **The kernel must be careful not to be tricked into dereferencing such a pointer, since that might reveal private information or destroy the integrity of the kernel.**

### How to solve this problem?

1. When a program passes the kernel a pointer, the kernel will check that the address is in the user part of the address space, and that the page table would allow the memory operation. Thus, the kernel will never suffer a page fault due to dereferencing a user-supplied pointer.

2. If the kernel does page fault, it should panic and terminate.

### 4.5.2   Exercise 9

**Exercise 9.** Change `kern/trap.c` to panic if a page fault happens in kernel mode.

Hint: to determine whether a fault happened in user mode or in kernel mode, check the low bits of the `tf_cs`.

Read `user_mem_assert` in `kern/pmap.c` and implement `user_mem_check` in that same file.

Change `kern/syscall.c` to sanity check arguments to system calls.

Boot your kernel, running `user/buggyhello`. The environment should be destroyed, and the kernel should *not* panic. You should see:

```
[00001000] user_mem_check assertion failure for va 00000001
[00001000] free env 00001000
Destroyed the only environment - nothing more to do!
```

Finally, change `debuginfo_eip` in `kern/kdebug.c` to call `user_mem_check` on `usd`, `stabs`, and `stabstr`. If you now run `user/breakpoint`, you should be able to run `backtrace` from the kernel monitor and see the backtrace traverse into `lib/libmain.c` before the kernel panics with a page fault. What causes this page fault? You don't need to fix it, but you should understand why it happens.

Figure 6: Exercise 9

We have the following steps:

1. Check it is in the kernel mode or user mode: check low bits of the **tf_cs**.

2. Implement **user_mem_assert()** in **kern/pmap.c**.

3. ddd

The following is about the code work:

1. **Check the mode:**

```
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          // We've already handled kernel-mode exceptions, so if we get here,
17          // the page fault happened in user mode.
18
19          // Destroy the environment that caused the fault.
20          cprintf("[%08x] user fault va %08x ip %08x\n",
21                  curenv->env_id, fault_va, tf->tf_eip);
22          print_trapframe(tf);
23          env_destroy(curenv);
24 }
```

Listing 2: page_fault_handler() in kern/trap.c

2. **Check the user_mem_assert()**

```
1  //
2  // Checks that environment 'env' is allowed to access the range
3  // of memory [va, va+len) with permissions 'perm | PTE_U | PTE_P'.
4  // If it can, then the function simply returns.
5  // If it cannot, 'env' is destroyed and, if env is the current
6  // environment, this function will not return.
7  //
8  void
9  user_mem_assert(struct Env *env, const void *va, size_t len, int perm)
10 {
11         if (user_mem_check(env, va, len, perm | PTE_U) < 0) {
12                 cprintf("[%08x] user_mem_check assertion failure for "
13                         "va %08x\n", env->env_id, user_mem_check_addr);
14                 env_destroy(env);       // may not return
15         }
16 }
```

It calls user_mem_check()

```
1  //
2  // Check that an environment is allowed to access the range of memory
3  // [va, va+len) with permissions 'perm | PTE_P'.
4  // Normally 'perm' will contain PTE_U at least, but this is not required.
5  // 'va' and 'len' need not be page-aligned; you must test every page that
6  // contains any of that range.  You will test either 'len/PGSIZE',
7  // 'len/PGSIZE + 1', or 'len/PGSIZE + 2' pages.
8  //
9  // A user program can access a virtual address if (1) the address is below
10 // ULIM, and (2) the page table gives it permission.  These are exactly
11 // the tests you should implement here.
12 //
13 // If there is an error, set the 'user_mem_check_addr' variable to the first
14 // erroneous virtual address.
```

```
15  //
16  // Returns 0 if the user program can access this range of addresses,
17  // and -E_FAULT otherwise.
18  //
19  int
20  user_mem_check(struct Env *env, const void *va, size_t len, int perm)
21  {
22          // LAB 3: Your code here.
23
24          return 0;
25  }
```

See its implementation:

```
1   int
2   user_mem_check(struct Env *env, const void *va, size_t len, int perm)
3   {
4           // LAB 3: Your code here.
5           char *start = NULL;
6           char *end = NULL;
7           start = ROUNDDOWN((char *)va, PGSIZE);
8           end = ROUNDUP((char *)(va + len), PGSIZE);
9
10          pte_t *cur_pg = NULL;
11          for(; start<end; start+=PGSIZE) {
12                  cur_pg = pgdir_walk(env->env_pgdir, (void *)start, 0);
13                  if((int) start > ULIM || cur_pg == NULL ||(((uint32_t)(*cur_pg)
        ↪ & perm) != perm)) {
14                          if(start == ROUNDDOWN((char *)va, PGSIZE)) {
15                                  user_mem_check_addr = (uintptr_t) va;
16                          }
17                          else {
18                                  user_mem_check_addr = (uintptr_t) start;
19                          }
20                          return -E_FAULT;
21                  }
22          }
23
24          return 0;
25  }
```

Now, we finally finish this lab: 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