# Report [Lab 3 User Environments]

I have completed the challenge to continue execution from the current location and to perform single stepping.

## Exercise 1

The way to allocate space or map a region becomes quite familiar after lab 2

```
1    envs = (struct Env *) boot_alloc(sizeof(struct Env) * NENV);
2    // ...
3    boot_map_region(kern_pgdir, UENVS, PTSIZE, PADDR(envs), PTE_U | PTE_P);
```

I also update `page_init` since `envs` becomes the last data structure allocated by `boot_alloc` now.

```
1        if (i >= EXTPHYSMEM / PGSIZE && (uintptr_t)KADDR(i * PGSIZE) <
     (uintptr_t)(envs + NENV))
2            continue;
```

## Exercise 2

This exercise asks us to finish several functions in `env.c`

### `env_init()`

This function initializes the `env` array, setting `status` and `id` for each elements. And `env_free_list` is also set with the same order as in the array.

```
1   void
2   env_init(void)
3   {
4       // Set up envs array
5       // LAB 3: Your code here.
6       env_free_list = &envs[0];
7       for (int i = 0; i < NENV; i++)
8       {
9           envs[i].env_status = ENV_FREE;
10          envs[i].env_id = 0;
11          if (i) envs[i - 1].env_link = &envs[i];
12      }
13
14      // Per-CPU part of the initialization
15      env_init_percpu();
16  }
```

## env_setup_vm()

This function allocates a page directory for a new environment and set up the kernel portion of its virtual memory.

Since kernel portion of all environments is identical, we use `kern_pgdir` as a template to do this.

```
1    // LAB 3: Your code here.
2    e->env_pgdir = page2kva(p), p->pp_ref++;
3    for (int i = PDX(UTOP); i < NPDENTRIES; i++)
4    {
5        e->env_pgdir[i] = kern_pgdir[i];
6    }
```

## region_alloc()

This is a helper function to allocate a bulk of physic space for some contiguous virtual space. To make it more convenient we allow the start position and size not to align with `PGSIZE`. But this introduce some corner cases that if the function is call with two adjacent region of virtual address and they're not aligned, we must make sure that we won't allocate another page frame for the same page twice!

```
1    static void
2    region_alloc(struct Env *e, void *va, size_t len)
3    {
4        // LAB 3: Your code here.
5        // (But only if you need it for load_icode.)
6        //
7        // Hint: It is easier to use region_alloc if the caller can pass
8        //    'va' and 'len' values that are not page-aligned.
9        //    You should round va down, and round (va + len) up.
10       //    (Watch out for corner-cases!)
11       int retval;
12       uintptr_t L = (uintptr_t)ROUNDDOWN(va, PGSIZE);
13       uintptr_t R = (uintptr_t)ROUNDUP(va + len, PGSIZE);
14       for (uintptr_t i = L; i < R; i += PGSIZE)
15       {
16           pte_t* pte = pgdir_walk(e->env_pgdir, (void *)i, false);
17           if (pte != NULL && (*pte & PTE_P))
18           {
19               continue;
20           }
21           struct PageInfo *p = page_alloc(0);
22           if (p == NULL)
23           {
24               panic("region_alloc: page_alloc failed");
25           }
26           retval = page_insert(e->env_pgdir, p, (void *)i, PTE_U | PTE_W);
27           if (retval)
28           {
29               panic("region_alloc: %e", retval);
30           }
31       }
```

```
32  }
```

# load_icode()

This function parses an ELF binary image and loads it into the user address space, just like what we do at booting. `region_alloc` helps a lot on allocating space for each section of the binary program and we load the user's page directory into `%cr3` so we can use `memcpy` and `memset` to initialize contents conveniently.

Different from booting where we directly jump to the entry point of the kernel, we store the entry point in the `TrapFrame` and allocate a page for the user's stack.

```
1       // LAB 3: Your code here.
2       lcr3(PADDR(e->env_pgdir));
3       if ((((struct Elf *)binary)->e_magic != ELF_MAGIC)
4       {
5           panic("unrecoginized binary format");
6       }
7       struct Proghdr *ph, *eph;
8       ph = (struct Proghdr *) (binary + ((struct Elf *)binary)->e_phoff);
9       eph = ph + ((struct Elf *)binary)->e_phnum;
10      for (; ph < eph; ph++) if (ph->p_type == ELF_PROG_LOAD)
11      {
12          region_alloc(e, (void *)ph->p_va, ph->p_memsz);
13          memcpy((void *)(ph->p_va), binary + ph->p_offset, ph->p_filesz);
14          memset((void *)(ph->p_va + ph->p_filesz), 0, ph->p_memsz - ph-
    >p_filesz);
15      }
16
17      // Now map one page for the program's initial stack
18      // at virtual address USTACKTOP - PGSIZE.
19
20      // LAB 3: Your code here.
21      region_alloc(e, (void *)(USTACKTOP - PGSIZE), PGSIZE);
22      e->env_tf.tf_eip = ((struct Elf *)binary)->e_entry;
23      lcr3(PADDR(kern_pgdir));
```

# env_create()

This function is used to create the first user environment using `env_alloc` and `load_icode`.

```
1   void
2   env_create(uint8_t *binary, enum EnvType type)
3   {
4       // LAB 3: Your code here.
5       struct Env *e;
6       int retval = env_alloc(&e, 0);
7       if (retval)
8       {
9           panic("env_create: %e", retval);
10      }
11      load_icode(e, binary);
12      e->env_type = type;
13  }
```

## env_run()

This function is used to switch between tasks. We have to alter the status of current environment and the environment to run appropriately and restore the new environment's content by loading its page directory and register values.

```
1       // LAB 3: Your code here.
2
3       if (curenv != NULL)
4       {
5           if (curenv->env_status == ENV_RUNNING)
6           {
7               curenv->env_status = ENV_RUNNABLE;
8           }
9           // seems we don't have to save the content yet.
10      }
11      curenv = e;
12      e->env_status = ENV_RUNNING;
13      e->env_runs++;
14      lcr3(PADDR(e->env_pgdir));
15      env_pop_tf(&e->env_tf);
```

# Exercise 3

This exercise and following text in lab description give us some knowledge about the way 80386 use to handle interrupts. The main data structure is IDT (interrupt descriptor table) so that each interrupt can use its interrupt vector assigned as an index to find the function to handle it.

# Exercise 4

We have to set up IDT in this exercise which is a somehow burdensome work.

Two macros `TRAPHANDLER` and `TRAPHANDLER_NOEC` are given in `trapentry.S` as some traps in 80386 push an error code on stack and others don't. So `TRAPHANDLER_NOEC` does much the same thing as `TRAPHANDLER` including pushing the interrupt vector into stack and jumping to `_alltraps` except pushing an extra `0` in the same position of the error code so the 2 kinds of traps finally look like the same.

So we have to read the manual to see whether each trap has an error code and use the right macro to make its entry point.

```
1   /*
2    * Lab 3: Your code here for generating entry points for the different
        traps.
3    */
4
5   TRAPHANDLER_NOEC(ENTRY_DIVIDE, T_DIVIDE)
6   TRAPHANDLER_NOEC(ENTRY_DEBUG, T_DEBUG)
7   TRAPHANDLER_NOEC(ENTRY_NMI, T_NMI)   // Exception Error Code: Not applicable.
8   TRAPHANDLER_NOEC(ENTRY_BRKPT, T_BRKPT)
9   TRAPHANDLER_NOEC(ENTRY_OFLOW, T_OFLOW)
10  TRAPHANDLER_NOEC(ENTRY_BOUND, T_BOUND)
11  TRAPHANDLER_NOEC(ENTRY_ILLOP, T_ILLOP)
12  TRAPHANDLER_NOEC(ENTRY_DEVICE, T_DEVICE)
13  TRAPHANDLER(ENTRY_DBLFLT, T_DBLFLT)
14  TRAPHANDLER(ENTRY_TSS, T_TSS)
15  TRAPHANDLER(ENTRY_SEGNP, T_SEGNP)
16  TRAPHANDLER(ENTRY_STACK, T_STACK)
17  TRAPHANDLER(ENTRY_GPFLT, T_GPFLT)
18  TRAPHANDLER(ENTRY_PGFLT, T_PGFLT)
19  TRAPHANDLER_NOEC(ENTRY_FPERR, T_FPERR)
20  TRAPHANDLER(ENTRY_ALIGN, T_ALIGN)
21  TRAPHANDLER_NOEC(ENTRY_MCHK, T_MCHK)
22  TRAPHANDLER_NOEC(ENTRY_SIMDERR, T_SIMDERR)
```

And at `_alltraps`, we push the rest of values onto stack to form the layout of the data structure `TrapFrame`. Then we push the address of it (`%esp`) as the argument and call `trap`

```
1   /*
2    * Lab 3: Your code here for _alltraps
3    */
4
5   _alltraps:
6       pushw $0;
7       pushw %ds;
8       pushw $0;
9       pushw %es;
10      pushal;
11      movw $GD_KD, %ax;
12      movw %ax, %ds;
13      movw %ax, %es;
14      pushl %esp;
15      call trap;
```

In `trap_init` we are about to set up the IDT (i.e. the entry points for each interrupt handler we defines in `trapentry.S`). Actually the entries in IDT are called gate descriptors in 80386's terminology and there is one bit to determine whether the exception is a trap or a fault (they differ in the returning point after handling of the exception). So we need to refer to the manual again and use `SETGATE` macro to fill in the entries accordingly.

```
1       // LAB 3: Your code here.
```

```
2        SETGATE(idt[T_DIVIDE], 0, GD_KT, &ENTRY_DIVIDE, 0);
3        SETGATE(idt[T_DEBUG], 1, GD_KT, &ENTRY_DEBUG, 0);
4        // debug exception is preferred to be handled by a task
5        // "Instruction address breakpoint conditions are faults,
6        //  while other debug conditions are traps."
7        // set it to be trap here for simplicity
8        SETGATE(idt[T_NMI], 0, GD_KT, &ENTRY_NMI, 0); // exception class is not
   applicable for NMI
9        SETGATE(idt[T_BRKPT], 1, GD_KT, &ENTRY_BRKPT, 3);
10       SETGATE(idt[T_OFLOW], 1, GD_KT, &ENTRY_OFLOW, 0);
11       SETGATE(idt[T_BOUND], 0, GD_KT, &ENTRY_BOUND, 0);
12       SETGATE(idt[T_ILLOP], 0, GD_KT, &ENTRY_ILLOP, 0);
13       SETGATE(idt[T_DEVICE], 0, GD_KT, &ENTRY_DEVICE, 0);
14       SETGATE(idt[T_DBLFLT], 0, GD_KT, &ENTRY_DBLFLT, 0); // abort
15       SETGATE(idt[T_TSS], 0, GD_KT, &ENTRY_TSS, 0);
16       SETGATE(idt[T_SEGNP], 0, GD_KT, &ENTRY_SEGNP, 0);
17       SETGATE(idt[T_STACK], 0, GD_KT, &ENTRY_STACK, 0);
18       SETGATE(idt[T_GPFLT], 0, GD_KT, &ENTRY_GPFLT, 0);
19       SETGATE(idt[T_PGFLT], 0, GD_KT, &ENTRY_PGFLT, 0);
20       SETGATE(idt[T_FPERR], 0, GD_KT, &ENTRY_FPERR, 0);
21       SETGATE(idt[T_ALIGN], 0, GD_KT, &ENTRY_ALIGN, 0);
22       SETGATE(idt[T_MCHK], 0, GD_KT, &ENTRY_MCHK, 0); // abort
23       SETGATE(idt[T_SIMDERR], 0, GD_KT, &ENTRY_SIMDERR, 0);
```

Here I set the DPL field of all traps to be 0 to prevent user software to use `int` to invoke them in an unintended way except `T_BRKPT` so that user can use `int 3` instruction to debug.

## Questions

### 1

The only difference between the entry points for each trap is that whether any extra 0 is pushed onto stack to take the position of the error code which is only generated in some kinds of traps. If there is only one kind of entry point, we won't have a uniform structure `Trapframe` to store all the information we concern.

### 2

The DPL field of each trap is set to 0, preventing user programs to use `int` to invoke any of the traps. So `int $14` is an invalid action for user program `softint` and therefore general protection fault is invoked as a result of privilege violation. If this action is allowed, it means that user programs have a convenient way to give some signals to kernel in an unintended way such as invoking a page fault with a fabricated virtual address and something evil might stem from this.

## Exercise 5

As the work in last exercise make the control flow of all interrupts flow to `trap_dispath` with a unified `struct Trapframe` distinguished by `tf->tf_trapno`, it's simple to leverage this to let our `page_fault_handler` to handle that specific fault.

```
1      // in trap_dispatch
2      if (tf->tf_trapno == T_PGFLT)
3      {
4          page_fault_handler(tf);
5          return;
6      }
```

# Exercise 6

In `trap_dispatch()`, we call `monitor()` with `tf` for breakpoint exceptions.

```
1      if (tf->tf_trapno == T_BRKPT)
2      {
3          monitor(tf);
4          return;
5      }
```

# Challenge

The content in Chapter 12 Debugging is really useful for completing this challenge.

To understand how to continue execution from the monitor, we have to understand the control flow of interrupt handling in JOS. The entry points in IDT leads all interrupts to `_alltraps` and then `trap()`. In `trap()`, `trap_dispatch()` is called to dispatch interrupts to their true handlers. If `trap_dispatch()` returns without panicking, the control flow returns to the current environment running.

So to continue the execution we simple return -1 to force monitor ends its loop to eventually return from `trap_dispatch()`

```
1  int
2  mon_continue(int argc, char **argv, struct Trapframe *tf)
3  {
4      return -1;
5  }
```

To implement the function of single stepping, we have to understand the debug feature of 80386. There are two kinds of debug interrupt vectors: `T_DEBUG` and `T_BRKPT`. While we are familiar with the latter now, `T_DEBUG` is important for us to implement single stepping function. The debug exception can be triggered in various situations, defined by the debug registers `DR0` ~ `DR7` and some other bits, including the TF bit in `%eflags`. If the TF bit is set, debug exception will be triggered after every instruction completion.

So to implement single stepping, we implement `mon_si` to set the TF bit in the `eflags` image in `tf`. And modify `mon_continue` at the same time to clear the TF bit, leaving the single-step mode.

```
 1   int
 2   mon_continue(int argc, char **argv, struct Trapframe *tf)
 3   {
 4       tf->tf_eflags &= ~FL_TF;
 5       return -1;
 6   }
 7
 8   int
 9   mon_si(int argc, char **argv, struct Trapframe *tf)
10   {
11       tf->tf_eflags |= FL_TF;
12       return -1;
13   }
```

And in `trap_dispatch()`, we invoke the monitor for both `T_BRKPT` and `T_DEBUG`.

```
 1       if (tf->tf_trapno == T_BRKPT || tf->tf_trapno == T_DEBUG)
 2       {
 3           monitor(tf);
 4           return;
 5       }
```

# Questions

## 3

Since the interrupt/trap gate descriptor in 80386 has a field to determine the privilege level required to use `int` instruction to invoke it, a general protection fault will be invoked instead of breakpoint exception if we didn't set the corresponding entry in IDT with appropriate privilege level i.e. 3.

```
 1       SETGATE(idt[T_BRKPT], 1, GD_KT, &ENTRY_BRKPT, 3);
```

If this was set up incorrectly, there will be a violation of protection based on privilege and therefore a general protection fault is triggered.

## 4

For some trap/interrupt, we only wish it to be triggered in particular situation such as page fault for the access of a page not in RAM. So these mechanisms are introduced to avoid inducing such traps/interrupts in an unintended way which will introduce ambiguity in the handling of such trap/interrupt.

# Exercise 7

First we need to set up the IDT for the system call interrupt, which is without error code and returns to the instruction next to the one invoked the trap.

```
1   // trapentry.S
2   TRAPHANDLER_NOEC(ENTRY_SYSCALL, T_SYSCALL)
3
4   // trap.c, trap_init()
5       SETGATE(idt[T_SYSCALL], 1, GD_KT, &ENTRY_SYSCALL, 3);
```

As `lib/syscall.c` stores the arguments in registers and expects returning value in `%eax`, we manipulate the registers saved in `Trapframe` to achieve this.

```
1       if (tf->tf_trapno == T_SYSCALL)
2       {
3           int retval = syscall(tf->tf_regs.reg_eax,
4                           tf->tf_regs.reg_edx, tf->tf_regs.reg_ecx, tf-
    >tf_regs.reg_ebx,
5                           tf->tf_regs.reg_edi, tf->tf_regs.reg_esi);
6           if (retval < 0)
7           {
8               panic("trap_dispatch: %e", retval);
9           }
10          tf->tf_regs.reg_eax = retval;
11          return;
12      }
```

Finally, we call the corresponding kernel functions to handle the system call according to `syscallno` in `syscall()`

```
1       switch (syscallno) {
2       case SYS_cputs:
3           sys_cputs((char *)a1, a2);
4           return 0;
5       case SYS_cgetc:
6           return sys_cgetc();
7       case SYS_getenvid:
8           return sys_getenvid();
9       case SYS_env_destroy:
10          return sys_env_destroy(a1);
11      default:
12          return -E_INVAL;
13      }
```

# Exercise 8

In this exercise we have to set `thisenv` to point to corresponding element in `envs` array during user environment setup. `lib/entry/S` has defined `envs` to point to `UENVS` we set up in the virtual address space so we can use it to access the kernel's `envs` array. And the index is able to be required by the system call `sys_getenvid()` and `ENVX` macro defined at `inc/env.h`

```
1       // set thisenv to point at our Env structure in envs[].
2       // LAB 3: Your code here.
3       thisenv = &envs[ENVX(sys_getenvid())];
```

## Exercise 9 & Exercise 10

We treat page fault in the kernel as a kernel BUG, so we panic at such situation

```
1      // LAB 3: Your code here.
2      if (!(tf->tf_cs & 3))
3      {
4          panic("page fault in the kernel");
5      }
```

Since `user_mem_check` require to fill in `user_mem_check_addr` the first virtual address invalid, I implement a help function called `user_mem_check_page` to check whether one single virtual address is valid by acquiring corresponding PTE by `pgdir_walk`

```
1   static bool user_mem_check_page(struct Env *env, uintptr_t va, int perm)
2   {
3       if (va >= ULIM) return false; // kernel space
4       pte_t *pte = pgdir_walk(env->env_pgdir, (void *)va, false);
5       if (pte == NULL || ~(*pte) & PTE_P)
6       {
7           return false; // PTE not present
8       }
9       return ((*pte) & perm) == perm;
10  }
11  int
12  user_mem_check(struct Env *env, const void *va, size_t len, int perm)
13  {
14      // LAB 3: Your code here.
15      uintptr_t L = (uintptr_t)va;
16      uintptr_t R = L + len;
17
18      // check the first page
19      if (!user_mem_check_page(env, L, perm))
20      {
21          user_mem_check_addr = L;
22          return -E_FAULT;
23      }
24
25      // check the rest pages
26      uintptr_t i = L % PGSIZE == 0 ? L + PGSIZE : ROUNDUP(L, PGSIZE);
27      while (i < R)
28      {
29          if (!user_mem_check_page(env, i, perm))
30          {
31              user_mem_check_addr = i;
32              return -E_FAULT;
33          }
34          i += PGSIZE;
35      }
36
37      return 0;
38  }
```

Now we can check whether the pointer arguments passed to system calls are valid by calling `user_mem_assert`

```
1      // LAB 3: Your code here.
2      user_mem_assert(curenv, s, len, 0);
```

Do the same check in `debuginfo_eip()`. Just be careful for the pointer stuff.

```
1          // Make sure this memory is valid.
2          // Return -1 if it is not.  Hint: Call user_mem_check.
3          // LAB 3: Your code here.
4          if (user_mem_check(curenv, usd, sizeof(struct UserStabData), 0))
5          {
6              return -1;
7          }
8
9          stabs = usd->stabs;
10         stab_end = usd->stab_end;
11         stabstr = usd->stabstr;
12         stabstr_end = usd->stabstr_end;
13
14         // Make sure the STABS and string table memory is valid.
15         // LAB 3: Your code here.
16         if (user_mem_check(curenv, stabs, (stab_end - stabs) * sizeof(struct
   Stab), 0))
17         {
18             return -1;
19         }
20         if (user_mem_check(curenv, stabstr, stabstr_end - stabstr, 0))
21         {
22             return -1;
23         }
```

If now we run `user/breakpoint` and call `backtrace` command in the monitor, a page fault will happen in the kernel at printing the arguments for `libmain`.

```
1  Stack backtrace:
2    ebp effff00  eip f0100bf5  args 00000001 effff28 f01c5000 f0106ab7
   f0106876
3          kern/monitor.c:149: monitor+353
4    ebp effff80  eip f0104557  args f01c5000 effffbc f014b178 00000092
   f011bfd8
5          kern/trap.c:196: trap+312
6    ebp effffb0  eip f010463f  args effffbc 00000000 00000000 eebfdff0
   effffdc
7          kern/syscall.c:69: syscall+0
8    ebp eebfdff0  eip 00800031  args 00000000 00000000Incoming TRAP frame at
   0xeffffe74
9  kernel panic at kern/trap.c:277: page fault in the kernel
```

This is because the stack of `libmain` is at the top of user stack with only two arguments for `user/breakpoint` so when `backtrace` tries to print 5 words above `%eip` it accesses the invalid page above `UTOP` which causes a page fault in the kernel.