# MIT 6.828: Operating System Engineering

Lab Report

Lab 5: File system, Spawn and Shell

CSNLP

csnlp16@126.com

https://csnlp.github.io/

# Contents

# 1 Introduction

## 1.1 Auto Merge Failed

Actually, there are some conflict happen when I try to auto merge.

```
1  Auto-merging lib/printfmt.c
2  Auto-merging kern/trap.c
3  CONFLICT (content): Merge conflict in kern/trap.c
4  Auto-merging kern/syscall.c
5  Auto-merging kern/sched.c
6  Auto-merging kern/pmap.c
7  Auto-merging kern/init.c
8  CONFLICT (content): Merge conflict in kern/init.c
9  Auto-merging kern/env.c
10 CONFLICT (content): Merge conflict in kern/env.c
11 Automatic merge failed; fix conflicts and then commit the result.
```

### 1.1.1 env_create() in kern/env.c

```
1  void
2  env_create(uint8_t *binary, enum EnvType type)
3  {
4          // LAB 3: Your code here.
5  <<<<<<< HEAD
6
7          // If this is the file server (type == ENV_TYPE_FS) give it I/O privileges.
8          // LAB 5: Your code here.
9  =======
10         struct Env *env;
11         int rc;
12         // step 1: allocates a new env with env_alloc
13         rc = env_alloc(&env, 0);
14         if(rc != 0)
15                 panic("env_create: env_alloc failed");
16         // step 2: loads the named elf binary with load_icode
17         load_icode(env, binary);
18         // step 3: set env's env_type
19         env->env_type = type;
20
21 >>>>>>> lab4
22 }
```

Modify as

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

```
13          load_icode(env, binary);
14          // step 3: set env's env_type
15          env->env_type = type;
16
17          // If this is the file server (type == ENV_TYPE_FS) give it I/O privileges.
18          // LAB 5: Your code here.
19  }
```

### 1.1.2   i386_init() in kern/init.c

```
1   void
2   i386_init(void)
3   {
4           // Initialize the console.
5           // Can't call cprintf until after we do this!
6           cons_init();
7
8           cprintf("6828 decimal is %o octal!\n", 6828);
9
10          // Lab 2 memory management initialization functions
11          mem_init();
12
13          // Lab 3 user environment initialization functions
14          env_init();
15          trap_init();
16
17          // Lab 4 multiprocessor initialization functions
18          mp_init();
19          lapic_init();
20
21          // Lab 4 multitasking initialization functions
22          pic_init();
23
24          // Acquire the big kernel lock before waking up APs
25          // Your code here:
26          lock_kernel();
27
28          // Starting non-boot CPUs
29          boot_aps();
30
31          // Start fs.
32          ENV_CREATE(fs_fs, ENV_TYPE_FS);
33  #if defined(TEST)
34          // Don't touch -- used by grading script!
35          ENV_CREATE(TEST, ENV_TYPE_USER);
36  #else
37          // Touch all you want.
38  <<<<<<< HEAD
39          ENV_CREATE(user_icode, ENV_TYPE_USER);
40  =======
41          //ENV_CREATE(user_primes, ENV_TYPE_USER);
42          ENV_CREATE(user_yield, ENV_TYPE_USER);
43          ENV_CREATE(user_yield, ENV_TYPE_USER);
44          ENV_CREATE(user_yield, ENV_TYPE_USER);
45  >>>>>>> lab4
46  #endif // TEST*
47
48          // Should not be necessary - drains keyboard because interrupt has given up.
```

4

```
49          kbd_intr();
50
51          // Schedule and run the first user environment!
52          sched_yield();
53  }
```

Modify as:

```
1   void
2   i386_init(void)
3   {
4           // Initialize the console.
5           // Can't call cprintf until after we do this!
6           cons_init();
7
8           cprintf("6828 decimal is %o octal!\n", 6828);
9
10          // Lab 2 memory management initialization functions
11          mem_init();
12
13          // Lab 3 user environment initialization functions
14          env_init();
15          trap_init();
16
17          // Lab 4 multiprocessor initialization functions
18          mp_init();
19          lapic_init();
20
21          // Lab 4 multitasking initialization functions
22          pic_init();
23
24          // Acquire the big kernel lock before waking up APs
25          // Your code here:
26          lock_kernel();
27
28          // Starting non-boot CPUs
29          boot_aps();
30
31          // Start fs.
32          ENV_CREATE(fs_fs, ENV_TYPE_FS);
33
34  #if defined(TEST)
35          // Don't touch -- used by grading script!
36          ENV_CREATE(TEST, ENV_TYPE_USER);
37  #else
38          // Touch all you want.
39          ENV_CREATE(user_icode, ENV_TYPE_USER);
40          //ENV_CREATE(user_primes, ENV_TYPE_USER);
41          ENV_CREATE(user_yield, ENV_TYPE_USER);
42          ENV_CREATE(user_yield, ENV_TYPE_USER);
43          ENV_CREATE(user_yield, ENV_TYPE_USER);
44  #endif // TEST*
45
46          // Should not be necessary - drains keyboard because interrupt has given up.
47          kbd_intr();
48
49          // Schedule and run the first user environment!
50          sched_yield();
51  }
```

### 1.1.3 trap_dispatch() in kern/trap.c

```c
static void
trap_dispatch(struct Trapframe *tf)
{
        // Handle spurious interrupts
        // The hardware sometimes raises these because of noise on the
        // IRQ line or other reasons. We don't care.
        if (tf->tf_trapno == IRQ_OFFSET + IRQ_SPURIOUS) {
                cprintf("Spurious interrupt on irq 7\n");
                print_trapframe(tf);
                return;
        }
        // Handle processor exceptions.
        // LAB 3: Your code here.
        switch(tf->tf_trapno) {
                case T_PGFLT:
                        page_fault_handler(tf);
                        return;
                case T_BRKPT:
                        monitor(tf);
                        return;
                case T_SYSCALL:
                        tf->tf_regs.reg_eax = syscall(tf->tf_regs.reg_eax,
                                tf->tf_regs.reg_edx,
                                tf->tf_regs.reg_ecx, tf->tf_regs.reg_ebx,
                                tf->tf_regs.reg_edi, tf->tf_regs.reg_esi);
                        return;
                // Handle clock interrupts. Don't forget to acknowledge the
                // interrupt using lapic_eoi() before calling the scheduler!
                // LAB 4: Your code here.
                case IRQ_OFFSET + IRQ_TIMER:
                        lapic_eoi();
                        sched_yield();
                        return;
                default:
                        // 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;
                        }
                        break;
        }
<<<<<<< HEAD
        // Handle keyboard and serial interrupts.
        // LAB 5: 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);
```

```
55          return;
56      }
57 =======
58 >>>>>>> lab4
59 }
```

Modify as

```
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
46
47     // Unexpected trap: The user process or the kernel has a bug.
48     print_trapframe(tf);
49     if (tf->tf_cs == GD_KT)
50         panic("unhandled trap in kernel");
51     else {
```

```
52          env_destroy(curenv);
53          return;
54      }
55
56      // Handle keyboard and serial interrupts.
57      // LAB 5: Your code here.
58 }
```

## 1.2   New Files

| | |
|---|---|
| `fs/fs.c` | Code that mainipulates the file system's on–disk structure. |
| `fs/bc.c` | A simple block cache built on top of our user–level page fault handling facility. |
| `fs/ide.c` | Minimal PIO–based (non–interrupt–driven) IDE driver code. |
| `fs/serv.c` | The file system server that interacts with client environments using file system IPCs. |
| `lib/fd.c` | Code that implements the general UNIX–like file descriptor interface. |
| `lib/file.c` | The driver for on–disk file type, implemented as a file system IPC client. |
| `lib/console.c` | The driver for console input/output file type. |
| `lib/spawn.c` | Code skeleton of the `spawn` library call. |

Figure 1: New files for Lab 5

# 2 File System Preliminaries

## 2.1 On-Disk File System Structure

Most UNIX file system divide available disk space into two main types of regions:

- inode regions.

- data regions.

### 2.1.1 Sectors and Blocks

- **sectors**: Most disks cannot perform reads and writes at byte granularity and instead perform reads and writes in units of sectors.

- **blocks:** File system actually allocate and use disk storage in units of blocks.

Be wary of the distinction between the two terms: {sector, block}. Sector size is a property of the disk hardware, whereas block size is an aspect of the OS using the disk. In this sense, a file system's block size must be a multiple of the sector size of the underlying disk.

### 2.1.2 Super

See its implementation:

```
struct Super {
        uint32_t s_magic;               // Magic number: FS_MAGIC
        uint32_t s_nblocks;             // Total number of blocks on disk
        struct File s_root;             // Root directory node
};
```

### 2.1.3 File Meta-data

This meta-data includes:

- file's name.

- file's size.

- file's type.

- pointers to the blocks comprising the file.

```
struct File {
        char f_name[MAXNAMELEN];        // filename
        off_t f_size;                   // file size in bytes
        uint32_t f_type;                // file type

        // Block pointers.
        // A block is allocated iff its value is != 0.
        uint32_t f_direct[NDIRECT];     // direct blocks
        uint32_t f_indirect;            // indirect block

```

```
11        // Pad out to 256 bytes; must do arithmetic in case we're compiling
12        // fsformat on a 64-bit machine.
13        uint8_t f_pad[256 - MAXNAMELEN - 8 - 4*NDIRECT - 4];
14 } __attribute__((packed));      // required only on some 64-bit machines
```

### 2.1.4   Directories versus Regular Files

A File structure in our file system can represent either a regular file or a directory; these two types of "files" are distinguished by the type field in the File structure. The file system manages regular files and directory-files in exactly the same way, except that it does not interpret the contents of the data blocks associated with regular files at all, whereas the file system interprets the contents of a directory-file as a series of File structures describing the files and subdirectories within the directory.

# 3 The File System

## 3.1 Disk Access

### 3.1.1 Exercise 1

**Exercise 1.** `i386_init` identifies the file system environment by passing the type `ENV_TYPE_FS` to your environment creation function, `env_create`. Modify `env_create` in `env.c`, so that it gives the file system environment I/O privilege, but never gives that privilege to any other environment.

Make sure you can start the file environment without causing a General Protection fault. You should pass the "fs i/o" test in **make grade**.

#### Question

1. Do you have to do anything else to ensure that this I/O privilege setting is saved and restored properly when you subsequently switch from one environment to another? Why?

Figure 2: Exercise 1

#### 3.1.1.1 See the background knowledge

The x86 processor uses the IOPL bits in the EFLAGS register to determine whether protected-mode code is allowed to perform special device I/O instructions such as the IN and OUT instructions. Since all of the IDE disk registers we need to access are located in the x86's I/O space rather than being memory-mapped, giving "I/O privilege" to the file system environment is the only thing we need to do in order to allow the file system to access these registers. In effect, the IOPL bits in the EFLAGS register provides the kernel with a simple "all-or-nothing" method of controlling whether user-mode code can access I/O space. In our case, we want the file system environment to be able to access I/O space, but we do not want any other environments to be able to access I/O space at all.

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

```

```
13        struct Env *env;
14        int rc;
15        // step 1: allocates a new env with env_alloc
16        rc = env_alloc(&env, 0);
17        if(rc != 0)
18                panic("env_create: env_alloc failed");
19        // step 2: loads the named elf binary with load_icode
20        load_icode(env, binary);
21        // step 3: set env's env_type
22        env->env_type = type;
23
24        // If this is the file server (type == ENV_TYPE_FS) give it I/O privileges.
25        // LAB 5: Your code here.
26        if (type == ENV_TYPE_FS)
27                env->env_tf.tf_eflags |= FL_IOPL_MASK;
28
29 }
```

| Exercise | Score |
|----------|-------|
| Exercise 1 | 25/150 |

## 3.2   The Block Cache

### 3.2.1   Exercise 2

**Exercise 2.** Implement the `bc_pgfault` and `flush_block` functions in `fs/bc.c`.
`bc_pgfault` is a page fault handler, just like the one your wrote in the previous lab
for copy–on–write fork, except that its job is to load pages in from the disk in
response to a page fault. When writing this, keep in mind that (1) `addr` may not be
aligned to a block boundary and (2) `ide_read` operates in sectors, not blocks.

The `flush_block` function should write a block out to disk *if necessary*. `flush_block`
shouldn't do anything if the block isn't even in the block cache (that is, the page
isn't mapped) or if it's not dirty. We will use the VM hardware to keep track of
whether a disk block has been modified since it was last read from or written to
disk. To see whether a block needs writing, we can just look to see if the `PTE_D`
"dirty" bit is set in the `uvpt` entry. (The `PTE_D` bit is set by the processor in
response to a write to that page; see 5.2.4.3 in chapter 5 of the 386 reference
manual.) After writing the block to disk, `flush_block` should clear the `PTE_D` bit
using `sys_page_map`.

Use **make grade** to test your code. Your code should pass "check_bc",
"check_super", and "check_bitmap".

Figure 3: Exercise 2

### 3.2.2 Implementation of bc_pgfault() in fs/bc.c

#### 3.2.2.1 See its hints

```c
// Fault any disk block that is read in to memory by
// loading it from disk.
static void
bc_pgfault(struct UTrapframe *utf)
{
        void *addr = (void *) utf->utf_fault_va;
        uint32_t blockno = ((uint32_t)addr - DISKMAP) / BLKSIZE;
        int r;

        // Check that the fault was within the block cache region
        if (addr < (void*)DISKMAP || addr >= (void*)(DISKMAP + DISKSIZE))
                panic("page fault in FS: eip %08x, va %08x, err %04x",
                        utf->utf_eip, addr, utf->utf_err);

        // Sanity check the block number.
        if (super && blockno >= super->s_nblocks)
                panic("reading non-existent block %08x\n", blockno);

        // Allocate a page in the disk map region, read the contents
        // of the block from the disk into that page.
        // Hint: first round addr to page boundary. fs/ide.c has code to read
        // the disk.
        //
        // LAB 5: you code here:
        // Clear the dirty bit for the disk block page since we just read the
        // block from disk
        if ((r = sys_page_map(0, addr, 0, addr, uvpt[PGNUM(addr)] & PTE_SYSCALL)) < 0)
                panic("in bc_pgfault, sys_page_map: %e", r);

        // Check that the block we read was allocated. (exercise for
        // the reader: why do we do this *after* reading the block
        // in?)
        if (bitmap && block_is_free(blockno))
                panic("reading free block %08x\n", blockno);
}
```

#### 3.2.2.2 See its implementation

```c
static void
bc_pgfault(struct UTrapframe *utf)
{
        void *addr = (void *) utf->utf_fault_va;
        uint32_t blockno = ((uint32_t)addr - DISKMAP) / BLKSIZE;
        int r;

        // Check that the fault was within the block cache region
        if (addr < (void*)DISKMAP || addr >= (void*)(DISKMAP + DISKSIZE))
                panic("page fault in FS: eip %08x, va %08x, err %04x",
                        utf->utf_eip, addr, utf->utf_err);

        // Sanity check the block number.
        if (super && blockno >= super->s_nblocks)
                panic("reading non-existent block %08x\n", blockno);

```

```
17         // Allocate a page in the disk map region, read the contents
18         // of the block from the disk into that page.
19         // Hint: first round addr to page boundary. fs/ide.c has code to read
20         // the disk.
21         //
22         // LAB 5: you code here:
23         // first round addr to page boundary
24         addr = (void *) ROUNDDOWN(addr, PGSIZE);
25         // Allocate a page at addr
26         if ((r = sys_page_alloc(0, addr, PTE_P|PTE_W|PTE_U)) < 0)
27                 panic("bc_pgfault failed: sys_page_alloc failed: %e", r);
28         // IDE read
29         if((r = ide_read(blockno*BLKSECTS, addr, BLKSECTS)) < 0)
30                 panic("bc_pgfault failed: ide failed: %e", r);
31
32
33         // Clear the dirty bit for the disk block page since we just read the
34         // block from disk
35         if ((r = sys_page_map(0, addr, 0, addr, uvpt[PGNUM(addr)] & PTE_SYSCALL)) < 0)
36                 panic("in bc_pgfault, sys_page_map: %e", r);
37
38         // Check that the block we read was allocated. (exercise for
39         // the reader: why do we do this *after* reading the block
40         // in?)
41         if (bitmap && block_is_free(blockno))
42                 panic("reading free block %08x\n", blockno);
43 }
```

### 3.2.2.3   ide_read()

It use **ide_read()**, let's see

```
1  int
2  ide_read(uint32_t secno, void *dst, size_t nsecs)
3  {
4         int r;
5
6         assert(nsecs <= 256);
7
8         ide_wait_ready(0);
9
10        outb(0x1F2, nsecs);
11        outb(0x1F3, secno & 0xFF);
12        outb(0x1F4, (secno >> 8) & 0xFF);
13        outb(0x1F5, (secno >> 16) & 0xFF);
14        outb(0x1F6, 0xE0 | ((diskno&1)<<4) | ((secno>>24)&0x0F));
15        outb(0x1F7, 0x20);      // CMD 0x20 means read sector
16
17        for (; nsecs > 0; nsecs--, dst += SECTSIZE) {
18                if ((r = ide_wait_ready(1)) < 0)
19                        return r;
20                insl(0x1F0, dst, SECTSIZE/4);
21        }
22
23        return 0;
24 }
```

### 3.2.3 Implementation of flush_block()

```c
// Flush the contents of the block containing VA out to disk if
// necessary, then clear the PTE_D bit using sys_page_map.
// If the block is not in the block cache or is not dirty, does
// nothing.
// Hint: Use va_is_mapped, va_is_dirty, and ide_write.
// Hint: Use the PTE_SYSCALL constant when calling sys_page_map.
// Hint: Don't forget to round addr down.
void
flush_block(void *addr)
{
        uint32_t blockno = ((uint32_t)addr - DISKMAP) / BLKSIZE;

        if (addr < (void*)DISKMAP || addr >= (void*)(DISKMAP + DISKSIZE))
                panic("flush_block of bad va %08x", addr);

        // LAB 5: Your code here.
        panic("flush_block not implemented");
}
```

### 3.2.4 Present Grade

```
internal FS tests [fs/test.c]: OK (1.3s)
  fs i/o: OK
  check_bc: OK
  check_super: OK
  check_bitmap: OK
```

Grade: 40/150.

## 3.3 The Block Bitmap

After **fs_init** sets the bitmap pointer, we can treat bitmap as a packed array of bits, one for each block on the disk. See for example, **block_is_free**, which simply checks whether a given block is marked free in the bitmap.

### 3.3.1 Exercise 3

**Exercise 3.** Use `free_block` as a model to implement `alloc_block` in `fs/fs.c`, which should find a free disk block in the bitmap, mark it used, and return the number of that block. When you allocate a block, you should immediately flush the changed bitmap block to disk with `flush_block`, to help file system consistency.

Use `make grade` to test your code. Your code should now pass "alloc_block".

Figure 4: Exercise 3

### 3.3.2 The Implementation of alloc_block()

#### 3.3.2.1 See its hints

15

```
1  // Search the bitmap for a free block and allocate it.  When you
2  // allocate a block, immediately flush the changed bitmap block
3  // to disk.
4  //
5  // Return block number allocated on success,
6  // -E_NO_DISK if we are out of blocks.
7  //
8  // Hint: use free_block as an example for manipulating the bitmap.
9  int
10 alloc_block(void)
11 {
12         // The bitmap consists of one or more blocks.  A single bitmap block
13         // contains the in-use bits for BLKBITSIZE blocks.  There are
14         // super->s_nblocks blocks in the disk altogether.
15
16         // LAB 5: Your code here.
17         panic("alloc_block not implemented");
18         return -E_NO_DISK;
19 }
```

### 3.3.2.2   See the hints from free_block()

```
1  // Mark a block free in the bitmap
2  void
3  free_block(uint32_t blockno)
4  {
5         // Blockno zero is the null pointer of block numbers.
6         if (blockno == 0)
7                 panic("attempt to free zero block");
8         bitmap[blockno/32] |= 1<<(blockno%32);
9  }
```

### 3.3.2.3   See its implementation

```
1  int
2  alloc_block(void)
3  {
4         // The bitmap consists of one or more blocks.  A single bitmap block
5         // contains the in-use bits for BLKBITSIZE blocks.  There are
6         // super->s_nblocks blocks in the disk altogether.
7
8         // LAB 5: Your code here.
9         size_t i;
10        for(i = 1; i < super->s_nblocks; i++) {
11                // find the free block
12                if(block_is_free(i)) {
13                        bitmap[i/32] ^= (1 << (i%32));
14                        flush_block(&bitmap[i%32]);
15                        return i;
16                }
17        }
18        //panic("alloc_block not implemented");
19        return -E_NO_DISK;
20 }
```

### 3.3.3   Present Grade

```
1  internal FS tests [fs/test.c]: OK (1.3s)
2    fs i/o: OK
3    check_bc: OK
4    check_super: OK
5    check_bitmap: OK
6    alloc_block: OK
```

Grade: 45/150.

## 3.4   File Operations

### 3.4.1   Exercise 4

**Exercise 4.** Implement `file_block_walk` and `file_get_block`. `file_block_walk` maps from a block offset within a file to the pointer for that block in the `struct File` or the indirect block, very much like what `pgdir_walk` did for page tables. `file_get_block` goes one step further and maps to the actual disk block, allocating a new one if necessary.

Use `make grade` to test your code. Your code should pass "file_open", "file_get_block", and "file_flush/file_truncated/file rewrite", and "testfile".

Figure 5: Exercise 4

### 3.4.2   Functions in fs/fs.c

- **check_super()**: check the file system super-block.

```
1  // Validate the file system super-block.
2  void
3  check_super(void) {}
```

- **block_is_free()**: check if the block is free or not.

```
1  // Check to see if the block bitmap indicates that block 'blockno' is free.
2  // Return 1 if the block is free, 0 if not.
3  bool
4  block_is_free(uint32_t blockno)
```

- **free_block()**: mark a block free in the bit map.

```
1  // Mark a block free in the bitmap
2  void
3  free_block(uint32_t blockno)
```

- **alloc_block()**: search the bitmap for a free block and allocate it.

```
1  int
2  alloc_block(void)
```

- **check_bitmap()**

```
1  void
2  check_bitmap(void)
```

- **fs_init()**: initialize the file system.

```
1  void
2  fs_init(void)
```

- **file_block_walk()**: find the disk block number slot for the 'filebno'th block in file 'f'.

```
1  static int
2  file_block_walk(struct File *f, uint32_t filebno, uint32_t **ppdiskbno, bool
   ↪ alloc)
```

- **file_get_block()**: set *blk to the address in memory where the filebno'th block of file 'f' would be mapped.

```
1  int
2  file_get_block(struct File *f, uint32_t filebno, char **blk)
```

- **dir_lookup()**: try to find a file named "name" in dir. If so, set *file to it.

```
1  static int
2  dir_lookup(struct File *dir, const char *name, struct File **file)
```

- **dir_alloc_file**: set *file to point at a free File structure in dir.

```
1  static int
2  dir_alloc_file(struct File *dir, struct File **file)
```

- **skip_slash()**: Skip over slashes.

```
1  static const char*
2  skip_slash(const char *p)
```

- **walk_path**:

```
1  // Evaluate a path name, starting at the root.
2  // On success, set *pf to the file we found
3  // and set *pdir to the directory the file is in.
4  // If we cannot find the file but find the directory
5  // it should be in, set *pdir and copy the final path
6  // element into lastelem.
7  static int
8  walk_path(const char *path, struct File **pdir, struct File **pf, char *lastelem
   ↪ )
```

- **file_create**:

```
1  // Create "path".  On success set *pf to point at the file and return 0.
2  // On error return < 0.
3  int
4  file_create(const char *path, struct File **pf)
```

18

- **file_open**:

```
// Open "path".  On success set *pf to point at the file and return 0.
// On error return < 0.
int
file_open(const char *path, struct File **pf)
```

- **file_read**

```
// Read count bytes from f into buf, starting from seek position
// offset.  This meant to mimic the standard pread function.
// Returns the number of bytes read, < 0 on error.
ssize_t
file_read(struct File *f, void *buf, size_t count, off_t offset)
```

- **file_write**

```
// Write count bytes from buf into f, starting at seek position
// offset.  This is meant to mimic the standard pwrite function.
// Extends the file if necessary.
// Returns the number of bytes written, < 0 on error.
int
file_write(struct File *f, const void *buf, size_t count, off_t offset)
```

- **file_free_block()**:

```
// Remove a block from file f.  If it's not there, just silently succeed.
// Returns 0 on success, < 0 on error.
static int
file_free_block(struct File *f, uint32_t filebno)
```

- **file_set_size()**: set the size of file f, truncating or extending as necessary.

```
// Set the size of file f, truncating or extending as necessary.
int
file_set_size(struct File *f, off_t newsize)
```

- **file_flush()**:

```
// Flush the contents and metadata of file f out to disk.
// Loop over all the blocks in file.
// Translate the file block number into a disk block number
// and then check whether that disk block is dirty.  If so, write it out.
void
file_flush(struct File *f)
```

### 3.4.3   The Implementation of file_block_walk() in fs/fs.c

### 3.4.3.1   See its hints

```
// Find the disk block number slot for the 'filebno'th block in file 'f'.
// Set '*ppdiskbno' to point to that slot.
// The slot will be one of the f->f_direct[] entries,
// or an entry in the indirect block.
// When 'alloc' is set, this function will allocate an indirect block
// if necessary.
```

```
7    //
8    // Returns:
9    //       0 on success (but note that *ppdiskbno might equal 0).
10   //       -E_NOT_FOUND if the function needed to allocate an indirect block, but
11   //               alloc was 0.
12   //       -E_NO_DISK if there's no space on the disk for an indirect block.
13   //       -E_INVAL if filebno is out of range (it's >= NDIRECT + NINDIRECT).
14   //
15   // Analogy: This is like pgdir_walk for files.
16   // Hint: Don't forget to clear any block you allocate.
17   static int
18   file_block_walk(struct File *f, uint32_t filebno, uint32_t **ppdiskbno, bool alloc)
19   {
20           // LAB 5: Your code here.
21           panic("file_block_walk not implemented");
22   }
```

### 3.4.3.2 See its implementation

```
1    static int
2    file_block_walk(struct File *f, uint32_t filebno, uint32_t **ppdiskbno, bool alloc)
3    {
4            // LAB 5: Your code here.
5            // CASE 1: filebno is out of range, i.e., filebno >= NDIRECT + NINDIRECT
6            if (filebno >= NDIRECT + NINDIRECT)
7                    return -E_INVAL;
8
9            // CASE 2: direct
10           if (filebno < NDIRECT) {
11                   // find the disk block number slot for the
12                   // 'filebno'th block in file 'f'. Set '*ppdiskbno'
13                   // to point to that slot.
14                   if(ppdiskbno)
15                           *ppdiskbno = &(f->f_direct[filebno]);
16                   return 0;
17           }
18
19           // CASE 3: indirect
20           /*
21           when code come to this place, it means that
22           NDIRECT <= blockno < NDIRECT + NINDIRECT
23           */
24           // alloc = 0: want to allocate a indirect
25           // alloc = 1: want to allocate an direct
26           int blockno = 0;
27           // if f->f_indirect == 0, set the f->f_indirect
28           if (f->f_indirect == 0) {
29                   // if the function needed to allocate an indirect block
30                   // but alloc was 0. Return -E_NOT_FOUND
31                   if (alloc == 0)
32                           return -E_NOT_FOUND;
33                   // if there is no space on the disk for an indirect block
34                   if ((blockno = alloc_block()) < 0)
35                           return -E_NO_DISK;
36
37                   memset(diskaddr(blockno), 0, BLKSIZE);
38                   flush_block(diskaddr(blockno));
39                   f->f_indirect = blockno;
```

```
40          }
41          if (ppdiskbno)
42                  *ppdiskbno = (uint32_t*)diskaddr(f->f_indirect) + filebno - NDIRECT;
43          return 0;
44          //panic("file_block_walk not implemented");
45 }
```

### 3.4.4 The implementation of file_get_block()

#### 3.4.4.1 See its hints

```
1  // Set *blk to the address in memory where the filebno'th
2  // block of file 'f' would be mapped.
3  //
4  // Returns 0 on success, < 0 on error.  Errors are:
5  //      -E_NO_DISK if a block needed to be allocated but the disk is full.
6  //      -E_INVAL if filebno is out of range.
7  //
8  // Hint: Use file_block_walk and alloc_block.
9  int
10 file_get_block(struct File *f, uint32_t filebno, char **blk)
11 {
12         // LAB 5: Your code here.
13          panic("file_get_block not implemented");
14 }
```

#### 3.4.4.2 See its implementation

```
1  int
2  file_get_block(struct File *f, uint32_t filebno, char **blk)
3  {
4         // LAB 5: Your code here.
5         uint32_t *ppdiskbno;
6         int r;
7         if ((r = file_block_walk(f, filebno, &ppdiskbno, true)) < 0)
8                 return r;
9         // if (*ppdiskbno == 0), point nothing, pose it.
10        if (*ppdiskbno == 0) {
11                if ((r = alloc_block()) < 0)
12                        return -E_NO_DISK;
13                *ppdiskbno = r;
14                memset(diskaddr(r), 0, BLKSIZE);
15                flush_block(diskaddr(r));
16        }
17        *blk = diskaddr(*ppdiskbno);
18        return 0;
19        //panic("file_get_block not implemented");
20 }
```

## 3.5 The File System Interface

### 3.5.1 Some Background Knowledge

Now that we have the necessary functionality within the file system environment itself, we must make it accessible to other environments that wish to use the file system. Since other environments can't directly call functions in the file system environment, we'll expose access to the file system environment via a remote procedure call, or **RPC**, abstraction,

built atop JOS's IPC mechanism. Graphically, here's what a call to the file system server (say, read) looks like:
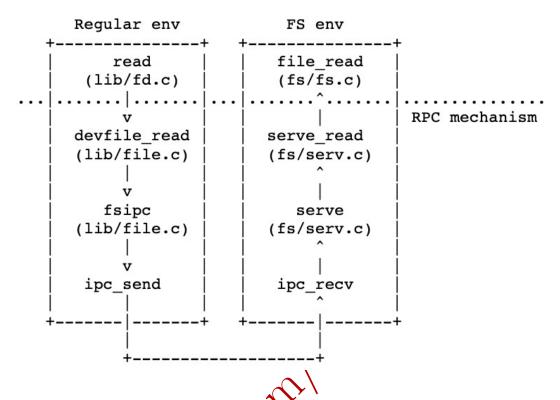
```
    Regular env              FS env
  +---------------+      +---------------+
  |     read      |      |   file_read   |
  |   (lib/fd.c)  |      |   (fs/fs.c)   |
...|.......|.......|...|.......^.......|..............
  |       v       |      |       |       |  RPC mechanism
  | devfile_read  |      |  serve_read   |
  | (lib/file.c)  |      |  (fs/serv.c)  |
  |       |       |      |       ^       |
  |       v       |      |       |       |
  |     fsipc     |      |     serve     |
  | (lib/file.c)  |      |  (fs/serv.c)  |
  |       |       |      |       ^       |
  |       v       |      |       |       |
  |   ipc_send    |      |   ipc_recv    |
  |       |       |      |       ^       |
  +-------|-------+      +-------|-------+
          |                      |
          +----------------------+
```

Figure 6: File System Server

See the whole process of call to the file system server:

- **read** works on any file descriptor and simply dispatches to the appropriate device read function, in this case, the **devfile_read** function.

- **decfile_read** implements **read** specifically for on-disk files.

- This and the other **devfile_*** functions in **lib/file.c** implment the client side of the FS operations and all work in roughly the same way, bundling up arguments in a request structure, calling **fsipc** to send the IPC request, and unpacking and returning the results.

### 3.5.2 Exercise 5

**Exercise 5.** Implement `serve_read` in `fs/serv.c`.

`serve_read`'s heavy lifting will be done by the already–implemented `file_read` in `fs/fs.c` (which, in turn, is just a bunch of calls to `file_get_block`). `serve_read` just has to provide the RPC interface for file reading. Look at the comments and code in `serve_set_size` to get a general idea of how the server functions should be structured.

Use **make grade** to test your code. Your code should pass "serve_open/file_stat/file_close" and "file_read" for a score of 70/150.

Figure 7: Exercise 5

### 3.5.3 The Implementation of serve_read in fs/serv.c

#### 3.5.3.1 See its hints

```
1  // Read at most ipc->read.req_n bytes from the current seek position
2  // in ipc->read.req_fileid.  Return the bytes read from the file to
3  // the caller in ipc->readRet, then update the seek position.  Returns
4  // the number of bytes successfully read, or < 0 on error.
5  int
6  serve_read(envid_t envid, union Fsipc *ipc)
7  {
8          struct Fsreq_read *req = &ipc->read;
9          struct Fsret_read *ret = &ipc->readRet;
10
11         if (debug)
12                 cprintf("serve_read %08x %08x %08x\n", envid, req->req_fileid, req->
      req_n);
13
14         // Lab 5: Your code here:
15         return 0;
16 }
```

The exercise5's hints suggest us to revisit the **serve_set_size()** code:

#### 3.5.3.2 serve_set_size() code

```
1  // Set the size of req->req_fileid to req->req_size bytes, truncating
2  // or extending the file as necessary.
3  int
4  serve_set_size(envid_t envid, struct Fsreq_set_size *req)
5  {
6          struct OpenFile *o;
7          int r;
8
9          if (debug)
```

23

```
10                    cprintf("serve_set_size %08x %08x %08x\n", envid, req->req_fileid, req
    ↪ ->req_size);
11
12            // Every file system IPC call has the same general structure.
13            // Here's how it goes.
14
15            // First, use openfile_lookup to find the relevant open file.
16            // On failure, return the error code to the client with ipc_send.
17            if ((r = openfile_lookup(envid, req->req_fileid, &o)) < 0)
18                    return r;
19
20            // Second, call the relevant file system function (from fs/fs.c).
21            // On failure, return the error code to the client.
22            return file_set_size(o->o_file, req->req_size);
23 }
```

Its steps:

1. Use **openfile_lookup** to find the relevant open file.

2. Call the relevant file system function.

### 3.5.4 Present Grade

```
1 internal FS tests [fs/test.c]: OK (1.1s)
2    fs i/o: OK
3    check_bc: OK
4    check_super: OK
5    check_bitmap: OK
6    alloc_block: OK
7    file_open: OK
8    file_get_block: OK
9    file_flush/file_truncate/file rewrite: OK
10 testfile: OK (1.0s)
11    serve_open/file_stat/file_close: OK
12    file_read: OK
```

Grade: 70/150.

### 3.5.5 Exercise 6

**Exercise 6.** Implement `serve_write` in `fs/serv.c` and `devfile_write` in `lib/file.c`.

Use `make grade` to test your code. Your code should pass "file_write", "file_read after file_write", "open", and "large file" for a score of 90/150.

Figure 8: Exercise 6

### 3.5.6   The Implementation of serve_write() in fs/serv.c

#### 3.5.6.1   Let's see its hints

```
// Write req->req_n bytes from req->req_buf to req_fileid, starting at
// the current seek position, and update the seek position
// accordingly.  Extend the file if necessary.  Returns the number of
// bytes written, or < 0 on error.
int
serve_write(envid_t envid, struct Fsreq_write *req)
{
        if (debug)
                cprintf("serve_write %08x %08x %08x\n", envid, req->req_fileid, req->
    req_n);

        // LAB 5: Your code here.
        panic("serve_write not implemented");
}
```

#### 3.5.6.2   See its implementation

```
int
serve_write(envid_t envid, struct Fsreq_write *req)
{
        if (debug)
                cprintf("serve_write %08x %08x %08x\n", envid, req->req_fileid, req->
    req_n);

        // LAB 5: Your code here.
        int r;
        int req_n;
        struct OpenFile *of;

        //
        if ((r = openfile_lookup(envid, req->req_fileid, &of)) < 0)
                return r;
        req_n = req->req_n > PGSIZE ? PGSIZE : req->req_n;

        // Write req_n bytes from req->req_buf to
        // req_fileid, starting at current seek position.
        if ((r = file_write(of->o_file, req->req_buf, req_n, of->o_fd->fd_offset)) <
    0)
                return r;
        // update the seek position.
        of->o_fd->fd_offset += r;
        return r;

        //panic("serve_write not implemented");
}
```

### 3.5.7   The Implementation of devfile_write() in lib/file.c

#### 3.5.7.1   Firstly, see its hints

```
static ssize_t
devfile_write(struct Fd *fd, const void *buf, size_t n)
{
        // Make an FSREQ_WRITE request to the file system server.  Be
```

25

```
5          // careful: fsipcbuf.write.req_buf is only so large, but
6          // remember that write is always allowed to write *fewer*
7          // bytes than requested.
8          // LAB 5: Your code here
9
10         panic("devfile_write not implemented");
11 }
```

### 3.5.7.2   See its implementation

```
1  static ssize_t
2  devfile_write(struct Fd *fd, const void *buf, size_t n)
3  {
4          // Make an FSREQ_WRITE request to the file system server.  Be
5          // careful: fsipcbuf.write.req_buf is only so large, but
6          // remember that write is always allowed to write *fewer*
7          // bytes than requested.
8          // LAB 5: Your code here
9          int r;
10         n = MIN(n, sizeof(fsipcbuf.write.req_buf));
11
12         fsipcbuf.write.req_fileid = fd->fd_file.id;
13         fsipcbuf.write.req_n = n;
14         memmove(fsipcbuf.write.req_buf, buf, n);
15         return fsipc(FSREQ_WRITE, NULL);
16         //panic("devfile_write not implemented");
17 }
```

### 3.5.8   Present Grade

```
1    fs i/o: OK
2    check_bc: OK
3    check_super: OK
4    check_bitmap: OK
5    alloc_block: OK
6    file_open: OK
7    file_get_block: OK
8    file_flush/file_truncate/file rewrite: OK
9  testfile: OK (1.3s)
10   serve_open/file_stat/file_close: OK
11   file_read: OK
12   file_write: OK
13   file_read after file_write: OK
14   open: OK
15   large file: OK
```

Grade: 90/150.

# 4 Spawning Processes

## 4.1 Spawn

### 4.1.1 Exercise 7

**Exercise 7.** spawn relies on the new syscall sys_env_set_trapframe to initialize the state of the newly created environment. Implement sys_env_set_trapframe in kern/syscall.c (don't forget to dispatch the new system call in syscall()).

Test your code by running the user/spawnhello program from kern/init.c, which will attempt to spawn /hello from the file system.

Use make grade to test your code.

Figure 9: Exercise 7

#### 4.1.1.1 Hints in sys_env_set_trapframe()

```
1  // Set envid's trap frame to 'tf'.
2  // tf is modified to make sure that user environments always run at code
3  // protection level 3 (CPL 3), interrupts enabled, and IOPL of 0.
4  //
5  // Returns 0 on success, < 0 on error.  Errors are:
6  //      -E_BAD_ENV if environment envid doesn't currently exist,
7  //              or the caller doesn't have permission to change envid.
8  static int
9  sys_env_set_trapframe(envid_t envid, struct Trapframe *tf)
10 {
11         // LAB 5: Your code here.
12         // Remember to check whether the user has supplied us with a good
13         // address!
14         panic("sys_env_set_trapframe not implemented");
15 }
```

#### 4.1.1.2 See its implementation

```
1  static int
2  sys_env_set_trapframe(envid_t envid, struct Trapframe *tf)
3  {
4         // LAB 5: Your code here.
5         // Remember to check whether the user has supplied us with a good
6         // address!
7         struct Env *env;
8         int r;
9
10        if ((r = envid2env(envid, &env, 1)) < 0)
```

```
11              return -E_BAD_ENV;
12      user_mem_assert(env, tf, sizeof(struct Trapframe), PTE_U);
13      env->env_tf = *tf;
14      env->env_tf.tf_cs |= 3;
15      env->env_tf.tf_eflags |= FL_IF;
16      return 0;
17      //panic("sys_env_set_trapframe not implemented");
18 }
```

### 4.1.2 DO NOT FORGET syscall in kern/syscall.c

```
1      // LAB 5: sys_env_set_trapframe
2      case SYS_env_set_trapframe:
3              return sys_env_set_trapframe((envid_t)a1, (struct Trapframe*)a2);
```

### 4.1.3 The Results

```
1  cui@cui-VirtualBox:~/mit6828/lab$ make run-spawnhello
2  sed "s/localhost:1234/localhost:26000/" < .gdbinit.tmpl > .gdbinit
3  make[1]: Entering directory '/home/cui/mit6828/lab'
4  + cc kern/init.c
5  + ld obj/kern/kernel
6  + mk obj/kern/kernel.img
7  make[1]: 'obj/fs/fs.img' is up to date.
8  make[1]: Leaving directory '/home/cui/mit6828/lab'
9  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 -drive file=obj/fs/fs.img,index
       ↪  =1,media=disk,format=raw
10 6828 decimal is 15254 octal!
11 Physical memory: 131072K available, base = 640K, extended = 130432K
12 check_page_free_list() succeeded!
13 check_page_alloc() succeeded!
14 check_page() succeeded!
15 check_kern_pgdir() succeeded!
16 check_page_free_list() succeeded!
17 check_page_installed_pgdir() succeeded!
18 SMP: CPU 0 found 1 CPU(s)
19 enabled interrupts: 1 2 4
20 FS is running
21 FS can do I/O
22 Device 1 presence: 1
23 i am parent environment 00001001
24 block cache is good
25 superblock is good
26 bitmap is good
27 alloc_block is good
28 file_open is good
29 file_get_block is good
30 file_flush is good
31 file_truncate is good
32 file rewrite is good
33 hello, world
34 i am environment 00001002
35 No runnable environments in the system!
36 Welcome to the JOS kernel monitor!
37 Type 'help' for a list of commands.
38 K>
```

```
1  internal FS tests [fs/test.c]: OK (1.7s)
2    fs i/o: OK
3    check_bc: OK
4    check_super: OK
5    check_bitmap: OK
6    alloc_block: OK
7    file_open: OK
8    file_get_block: OK
9    file_flush/file_truncate/file rewrite: OK
10 testfile: OK (1.2s)
11   serve_open/file_stat/file_close: OK
12   file_read: OK
13   file_write: OK
14   file_read after file_write: OK
15   open: OK
16   large file: OK
17 spawn via spawnhello: OK (0.9s)
18     (Old jos.out.spawn failure log removed)
```

Grade: 95/150.

## 4.2 Sharing Library State across Fork and Spawn

### 4.2.1 Background Knowledge

In JOS, each of these device types has a corresponding struct Dev

### 4.2.2 Exercise 8

**Exercise 8.** Change `duppage` in `lib/fork.c` to follow the new convention. If the page table entry has the `PTE_SHARE` bit set, just copy the mapping directly. (You should use `PTE_SYSCALL`, not `0xfff`, to mask out the relevant bits from the page table entry. `0xfff` picks up the accessed and dirty bits as well.)

Likewise, implement `copy_shared_pages` in `lib/spawn.c`. It should loop through all page table entries in the current process (just like `fork` did), copying any page mappings that have the `PTE_SHARE` bit set into the child process.

Figure 10: Exercise 8

### 4.2.3 The Implementation of duppage() in lib/fork.c: PTE_SHARE VERSION

#### 4.2.3.1 See the PTE_SHARE version

```
1 //
2 // Map our virtual page pn (address pn*PGSIZE) into the target envid
3 // at the same virtual address.  If the page is writable or copy-on-write,
```

```
 4  // the new mapping must be created copy-on-write, and then our mapping must be
 5  // marked copy-on-write as well.  (Exercise: Why do we need to mark ours
 6  // copy-on-write again if it was already copy-on-write at the beginning of
 7  // this function?)
 8  //
 9  // Returns: 0 on success, < 0 on error.
10  // It is also OK to panic on error.
11  //
12  static int
13  duppage(envid_t envid, unsigned pn)
14  {
15          int r;
16          // LAB 4: Your code here.
17  //
18  //      void *addr;
19  //      pte_t pte;
20  //      int perm;
21  //
22  //      addr = (void *)((uint32_t)pn * PGSIZE);
23  //      pte = uvpt[pn];
24  //      perm = PTE_P | PTE_U;
25  //      if ((pte & PTE_W) || (pte & PTE_COW))
26  //              perm |= PTE_COW;
27  //      if ((r = sys_page_map(thisenv->env_id, addr, envid, addr, perm)) < 0) {
28  //              panic("duppage failed: page remapping failed: %e", r);
29  //              return r;
30  //      }
31  //      if (perm & PTE_COW) {
32  //              if ((r = sys_page_map(thisenv->env_id, addr, thisenv->env_id, addr,
        ↪ perm)) < 0) {
33  //                      panic("duppage failed: page remapping failed: %e", r);
34  //                      return r;
35  //              }
36  //      }
37  //      //panic("duppage not implemented");
38  //      return 0;
39          // LAB 5: For exercise 8.
40          void *addr;
41          pte_t pte;
42          int perm;
43
44          addr = (void *)((uint32_t)pn * PGSIZE);
45          pte = uvpt[pn];
46          // If the page table entry has the PTE_SHARE bit set,
47          // just copy the mapping directly.
48          if (pte & PTE_SHARE) {
49                  if ((r = sys_page_map(sys_getenvid(), addr, envid, addr, pte &
        ↪ PTE_SYSCALL)) < 0) {
50                          panic("duppage failed: sys_page_map failed: %e", r);
51                          return r;
52                  }
53          }
54          else {
55                  perm = PTE_P | PTE_U;
56                  if ((pte & PTE_W) || (pte & PTE_COW))
57                          perm |= PTE_COW;
58                  if ((r = sys_page_map(thisenv->env_id, addr, envid, addr, perm)) < 0)
        ↪ {
```

```
59                        panic("duppage failed: page remapping failed: %e", r);
60                        return r;
61                }
62                if (perm & PTE_COW) {
63                        if ((r = sys_page_map(thisenv->env_id, addr, thisenv->env_id,
     addr, perm)) < 0) {
64                                panic("duppage failed: page remapping failed: %e", r);
65                                return r;
66                        }
67                }
68        }
69        return 0;
70 }
```

## 4.2.4   The Implementation of copy_shared_pages() in lib/spawn.c: PTE_SHARE VERSION

It loops through all page table entries in the current process (just like the **fork** did), copying any page mappings that have the **PTE_SHARE** bit set into the child process.

### 4.2.4.1   See its hints

# 5 The Keyboard Interface

## 5.1 Exercise 9

**Exercise 9.** In your `kern/trap.c`, call `kbd_intr` to handle trap `IRQ_OFFSET+IRQ_KBD` and `serial_intr` to handle trap `IRQ_OFFSET+IRQ_SERIAL`.

Figure 11: Exercise 9

See the **trap_dispatch()**

```
// Handle keyboard and serial interrupts.
// LAB 5: Your code here.
case (IRQ_OFFSET + IRQ_KBD):
        lapic_eoi();
        kbd_intr();
        return;
case (IRQ_OFFSET + IRQ_SERIAL):
        lapic_eoi();
        serial_intr();
        return;
```

# 6 The Shell

## 6.1 runcmd in user/sh.c

```
1  // Parse a shell command from string 's' and execute it.
2  // Do not return until the shell command is finished.
3  // runcmd() is called in a forked child,
4  // so it's OK to manipulate file descriptor state.
5  #define MAXARGS 16
6  void
7  runcmd(char* s)
8  {
9          char *argv[MAXARGS], *t, argv0buf[BUFSIZ];
10         int argc, c, i, r, p[2], fd, pipe_child;
11
12         pipe_child = 0;
13         gettoken(s, 0);
14
15  again:
16         argc = 0;
17         while (1) {
18                 switch ((c = gettoken(0, &t))) {
19
20                 case 'w':       // Add an argument
21                     if (argc == MAXARGS) {
22                             cprintf("too many arguments\n");
23                             exit();
24                     }
25                     argv[argc++] = t;
26                     break;
27                 case '<':       // Input redirection
28                     // Grab the filename from the argument list
29                     if (gettoken(0, &t) != 'w') {
30                             cprintf("syntax error: < not followed by word\n");
31                             exit();
32                     }
33                     // Open 't' for reading as file descriptor 0
34                     // (which environments use as standard input).
35                     // We can't open a file onto a particular descriptor,
36                     // so open the file as 'fd',
37                     // then check whether 'fd' is 0.
38                     // If not, dup 'fd' onto file descriptor 0,
39                     // then close the original 'fd'.
40
41                     // LAB 5: Your code here.
42                     panic("< redirection not implemented");
43                     break;
44                 case '>':       // Output redirection
45                     // Grab the filename from the argument list
46                     if (gettoken(0, &t) != 'w') {
47                             cprintf("syntax error: > not followed by word\n");
48                             exit();
49                     }
50                     if ((fd = open(t, O_WRONLY|O_CREAT|O_TRUNC)) < 0) {
51                             cprintf("open %s for write: %e", t, fd);
52                             exit();
53                     }
```

```
54                    if (fd != 1) {
55                            dup(fd, 1);
56                            close(fd);
57                    }
58                    break;
59            case '|':        // Pipe
60                    if ((r = pipe(p)) < 0) {
61                            cprintf("pipe: %e", r);
62                            exit();
63                    }
64                    if (debug)
65                            cprintf("PIPE: %d %d\n", p[0], p[1]);
66                    if ((r = fork()) < 0) {
67                            cprintf("fork: %e", r);
68                            exit();
69                    }
70                    if (r == 0) {
71                            if (p[0] != 0) {
72                                    dup(p[0], 0);
73                                    close(p[0]);
74                            }
75                            close(p[1]);
76                            goto again;
77                    } else {
78                            pipe_child = r;
79                            if (p[1] != 1) {
80                                    dup(p[1], 1);
81                                    close(p[1]);
82                            }
83                            close(p[0]);
84                            goto runit;
85                    }
86                    panic("| not implemented");
87                    break;

89            case 0:          // String is complete
90                    // Run the current command!
91                    goto runit;

93            default:
94                    panic("bad return %d from gettoken", c);
95                    break;

97            }
98        }

100 runit:
101        // Return immediately if command line was empty.
102        if(argc == 0) {
103                if (debug)
104                        cprintf("EMPTY COMMAND\n");
105                return;
106        }

108        // Clean up command line.
109        // Read all commands from the filesystem: add an initial '/' to
110        // the command name.
111        // This essentially acts like 'PATH=/'.
```

34

```
112        if (argv[0][0] != '/') {
113                argv0buf[0] = '/';
114                strcpy(argv0buf + 1, argv[0]);
115                argv[0] = argv0buf;
116        }
117        argv[argc] = 0;
118
119        // Print the command.
120        if (debug) {
121                cprintf("[%08x] SPAWN:", thisenv->env_id);
122                for (i = 0; argv[i]; i++)
123                        cprintf(" %s", argv[i]);
124                cprintf("\n");
125        }
126        // Spawn the command!
127        if ((r = spawn(argv[0], (const char**) argv)) < 0)
128                cprintf("spawn %s: %e\n", argv[0], r);
129
130        // In the parent, close all file descriptors and wait for the
131        // spawned command to exit.
132        close_all();
133        if (r >= 0) {
134                if (debug)
135                        cprintf("[%08x] WAIT %s %08x\n", thisenv->env_id, argv[0], r);
136                wait(r);
137                if (debug)
138                        cprintf("[%08x] wait finished\n", thisenv->env_id);
139        }
140
141        // If we were the left-hand part of a pipe,
142        // wait for the right-hand part to finish.
143        if (pipe_child) {
144                if (debug)
145                        cprintf("[%08x] WAIT pipe_child %08x\n", thisenv->env_id,
     ↪ pipe_child);
146                wait(pipe_child);
147                if (debug)
148                        cprintf("[%08x] wait finished\n", thisenv->env_id);
149        }
150
151        // Done!
152        exit();
153 }
```

See the implementation:

```
1                    // Open 't' for reading as file descriptor 0
2                    // (which environments use as standard input).
3                    // We can't open a file onto a particular descriptor,
4                    // so open the file as 'fd',
5                    // then check whether 'fd' is 0.
6                    // If not, dup 'fd' onto file descriptor 0,
7                    // then close the original 'fd'.
8
9                    // LAB 5: Your code here.
10                   if ((fd = open(t, O_RDONLY)) < 0) {
11                           cprintf("open file: %s failed: %e", t, fd);
12                           exit();
13                   }
```

```
14                          if (fd != 0) {
15                                  dup(fd, 0);
16                                  close(fd);
17                          }
18                          //panic("< redirection not implemented");
19                          break;
```

## 6.2 Big BUG

### 6.2.1 Present Grade

```
1  make[1]: Leaving directory '/home/cui/mit6828/lab'
2  internal FS tests [fs/test.c]: OK (2.1s)
3    fs i/o: OK
4    check_bc: OK
5    check_super: OK
6    check_bitmap: OK
7    alloc_block: OK
8    file_open: OK
9    file_get_block: OK
10   file_flush/file_truncate/file rewrite: OK
11 testfile: OK (1.3s)
12   serve_open/file_stat/file_close: OK
13   file_read: OK
14   file_write: OK
15   file_read after file_write: OK
16   open: OK
17   large file: OK
18 spawn via spawnhello: OK (1.1s)
19 Protection I/O space: FAIL (1.3s)
20     AssertionError: ...
21         (null): made it here --- bug
22         No runnable environments in the system!
23         Welcome to the JOS kernel monitor!
24         Type 'help' for a list of commands.
25         qemu: terminating on signal 15 from pid 26196
26     MISSING 'TRAP'
27
28     QEMU output saved to jos.out.faultio
29 PTE_SHARE [testpteshare]: OK (1.2s)
30 PTE_SHARE [testfdsharing]: OK (1.9s)
31 start the shell [icode]: Timeout! OK (30.7s)
32 testshell: OK (3.2s)
33 primespipe: OK (8.9s)
34 Score: 145/150
35 make: *** [grade] Error 1
```

# 7 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

https://github.com/csnlp/MIT6828