

Report [Lab 5: File system, Spawn and Shell]

I've completed the challenge of `mmap`

The File System

Disk Access

Exercise 1

An environment is permitted to do IO only when its CPL is equal or less than the IOPL field in EFLAGS[1]. As the file system environment is running in user mode with CPL=3, we should set the IOPL to 3 to give it appropriate privilege.

```
1  if (type == ENV_TYPE_FS)
2  {
3      e->env_tf.tf_eflags |= FL_IOPL_3;
4  }
```

Question 1

Nothing special needs to be done since IOPL is part of EFLAGS which is properly stored in `struct Trapframe` when trapped into kernel and restored when to be run again.

The Block Cache

Exercise 2

In JOS we map the whole disk into the address space of the file system environment and read a page (which has the same size as a block) from disk when necessary i.e. it is accessed by leveraging the `pgfault_handler` mechanism.

We use `sys_page_alloc` to allocate a new page for the faulting address and use `ide_read` to load its content from the disk. As disk hardware is working with the unit of sector, we pass `BLKSECTS` as the `nsecs` parameter to read the sectors of the corresponding block.

```
1  // LAB 5: you code here:
2  addr = ROUNDDOWN(addr, PGSIZE);
3  if ((r = sys_page_alloc(0, addr, PTE_U | PTE_P | PTE_W)) < 0)
4  {
5      panic("in bc_pgfault, sys_page_alloc: %e", r);
6  }
7  if ((r = ide_read(blockno * BLKSECTS, addr, BLKSECTS)) < 0)
8  {
9      panic("in bc_pgfault, ide_read: %e", r);
10 }
```

`flush_block()` will write the content cached in memory back to disk if necessary, that is to say, the corresponding page is mapped and modified (dirty). `ide_write()` with similar semantics as `ide_read` is used to accomplish the actual work. The dirty bit also needs to be cleared when things are done.

```
1 // LAB 5: Your code here.
2 int r;
3 addr = ROUNDDOWN(addr, PGSIZE);
4 if (va_is_mapped(addr) && va_is_dirty(addr))
5 {
6     if ((r = ide_write(blockno * BLKSECTS, addr, BLKSECTS)) < 0)
7     {
8         panic("in flush_block, ide_write: %e", r);
9     }
10    // clear dirty bit by clearing all bits in PTE except permission
    bits
11    if ((r = sys_page_map(0, addr, 0, addr, uvpt[PGNUM(addr)] &
    PTE_SYSCALL)) < 0)
12    {
13        panic("in flush_block, sys_page_map: %e", r);
14    }
15 }
```

The Block Bitmap

Exercise 3

Bitmap, as a critical data structure for the disk, should be flushed back immediately after any modification. Keeping this in mind, we iterate over `super->s_nblocks` block numbers and use `block_is_free` to search for a free block. Once it's found, we revert the corresponding bit in the bitmap to mark its allocation and flush this change to disk to maintain consistency of the disk.

```
1 // LAB 5: Your code here.
2 // panic("alloc_block not implemented");
3 if (super == 0)
4 {
5     panic("alloc_block: super is null");
6 }
7 for (int i = 0; i < super->s_nblocks; i++)
8 {
9     if (block_is_free(i))
10    {
11        bitmap[i / 32] ^= 1 << (i % 32);
12        flush_block(&bitmap[i / 32]);
13        return i;
14    }
15 }
16 return -E_NO_DISK;
```

File Operations

Exercise 4

File system provides the abstraction of files to applications to make disk management more convenient. To support this abstraction, we must have corresponding mechanism to translation a position in the file to a position in the disk just like what we do to translate virtual address to physical address.

The underlying functions are `file_block_walk()` and `file_get_block()`. In `file_block_walk()`, a `filebno` is given and the pointer to where the block number of `filebno`-th block in corresponding file is stored in the file system is to be filled. If `filebno < NDIRECT`, the position of the block is stored in `struct File`. Otherwise the block number is stored in the indirect block which may need to be allocated depending on its existence and `alloc` parameter.

```
1 // Hint: Don't forget to clear any block you allocate.
2 static int
3 file_block_walk(struct File *f, uint32_t filebno, uint32_t **ppdiskbno, bool
  alloc)
4 {
5     // LAB 5: Your code here.
6     // panic("file_block_walk not implemented");
7     int r;
8     if (filebno >= NDIRECT + NINDIRECT)
9     {
10         return -E_INVAL;
11     }
12     if (filebno < NDIRECT)
13     {
14         *ppdiskbno = f->f_direct + filebno;
15         return 0;
16     }
17     if (!f->f_indirect)
18     {
19         if (!alloc) return -E_NOT_FOUND;
20         if ((r = alloc_block()) < 0)
21         {
22             return r;
23         }
24         memset(diskaddr(r), 0, BLKSIZE);
25         f->f_indirect = r;
26     }
27     uint32_t *ind = diskaddr(f->f_indirect);
28     *ppdiskbno = ind + (filebno - NDIRECT);
29     return 0;
30 }
```

`file_get_block()` is basically a wrapper for `file_block_walk`, except that a new block must be allocated if the corresponding file block number is 0 i.e. null.

```
1 // Hint: Use file_block_walk and alloc_block.
2 int
3 file_get_block(struct File *f, uint32_t filebno, char **blk)
```

```

4  {
5      // LAB 5: Your code here.
6      // panic("file_get_block not implemented");
7      uint32_t *pdiskbno;
8      int r;
9      if ((r = file_block_walk(f, filebno, &pdiskbno, true)) < 0)
10     {
11         return r;
12     }
13     if (*pdiskbno == 0)
14     {
15         if ((r = alloc_block()) < 0)
16         {
17             return r;
18         }
19         memset(diskaddr(r), 0, sizeof(r));
20         *pdiskbno = r;
21     }
22     *blk = diskaddr(*pdiskbno);
23     return 0;
24 }

```

The file system interface

Exercise 5 & 6

In JOS, regular environments do file operations by RPC to the file system environment, which is implemented by IPC. A number is passed through IPC as an identifier to some specific file system service and the arguments is passed as a sharing page to file system environment, which is also used to contain the results.

In the server side, the arguments should be dispatched to the corresponding file system function wrapper and be parsed and passed to the function to get the result. Some other information needed is stored in the file descriptor field such as current offset for read and write operation.

```

1      // Lab 5: Your code here:
2      struct OpenFile *o;
3      int r;
4
5      if ((r = openfile_lookup(envid, req->req_fileid, &o)) < 0)
6          return r;
7
8      if ((r = file_read(o->o_file, ret->ret_buf, MIN(req->req_n, PGSIZE),
9                  o->o_fd->fd_offset)) < 0)
10         return r;
11
12     o->o_fd->fd_offset += r;
13     return r;

```

```

1 // LAB 5: Your code here.
2 // panic("serve_write not implemented");
3 struct OpenFile *o;
4 int r;
5
6 if ((r = openfile_lookup(envid, req->req_fileid, &o)) < 0)
7     return r;
8
9 if ((r = file_write(o->o_file, req->req_buf, req->req_n, o->o_fd-
>fd_offset)) < 0)
10     return r;
11
12 o->o_fd->fd_offset += r;
13 return r;

```

In the client side of RPC mechanism, arguments are written into the IPC page before passing them to server side. Then the client blocks until the file system environment finish the job and inform the client. Then some checks of returning value in the shared page is checked.

```

1 int
2 serve_write(envid_t envid, struct Fsreq_write *req)
3 {
4     if (debug)
5         cprintf("serve_write %08x %08x %08x\n", envid, req->req_fileid, req-
>req_n);
6
7     // LAB 5: Your code here.
8     // panic("serve_write not implemented");
9     struct OpenFile *o;
10    int r;
11
12    if ((r = openfile_lookup(envid, req->req_fileid, &o)) < 0)
13        return r;
14
15    if ((r = file_write(o->o_file, req->req_buf, req->req_n, o->o_fd-
>fd_offset)) < 0)
16        return r;
17
18    o->o_fd->fd_offset += r;
19    return r;
20 }

```

Spawning Processes

If we want to implement `exec` in user space, we have to ask the environment to replace itself including code and stacks to the new program's, which is difficult, if not impossible, to keep consistency in such procedure.

So instead we implement `spawn`, creating an empty child environment and replacing its content to the new program's, which is much easier.

Exercise 7

A new system call `sys_env_set_trapframe()` is needed by `spawn()` to fill the `env_tf` for the new environment to indicate the position of program entry point and initial stack. Some particular bits of the `struct Trapframe` is modified to ensure that the permission is appropriate for a user program.

```
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     // panic("sys_env_set_trapframe not implemented");
8
9     struct Env *e;
10    int r;
11
12    if ((r = envid2env(envid, &e, true)) < 0)
13        return r;
14
15    user_mem_assert(curenv, tf, sizeof(struct Trapframe), 0);
16    e->env_tf = *tf;
17
18    e->env_tf.tf_cs |= 3;
19    e->env_tf.tf_eflags |= FL_IF;
20    e->env_tf.tf_eflags &= ~FL_IOPL_MASK;
21
22    return 0;
23 }
```

challenge!

The `mmap` method defined in POSIX allows to map a file into a region in address space to make files able to be accessed like other data in memory. For read-only parts of files like the code in program image and archives, using `mmap` can make different environments share their common content to save space.

In JOS, I choose to implement `mmap` in a simple but effective way. It takes an offset and an address space which are both aligned, then it leverages file system call to share exactly one page through IPC, which is exactly the corresponding block in the mapped disk.

```
1 // fs/server.c
2 int
3 serve_read_map(envid_t envid, struct Fsreq_read_map *req,
4     void **pg_store, int *perm_store)
5 {
6     struct OpenFile *o;
7     char *blk;
8     int r;
9
10    if (req->req_perm & PTE_W)
11        return -E_INVALID;
```

```

12     if (req->req_offset % BLKSIZE)
13         return -E_INVAL;
14     int blkno = req->req_offset / BLKSIZE;
15
16     if ((r = openfile_lookup(envid, req->req_fileid, &o)) < 0)
17         return r;
18     if ((r = file_get_block(o->o_file, blkno, &blk)) < 0)
19         return r;
20     memcpy(blk, blk, PGSIZE); // load the block from the disk
21
22     *pg_store = blk;
23     *perm_store = req->req_perm;
24
25     return 0;
26 }

```

To use RPC mechanism to call this function, we have to register it in the user library, adding a new function pointer to `struct Dev` and implement client-side `read_map()` for files.

```

1 // inc/fd.h
2 // Per-device-class file descriptor operations
3 struct Dev {
4     int dev_id;
5     const char *dev_name;
6     ssize_t (*dev_read)(struct Fd *fd, void *buf, size_t len);
7     ssize_t (*dev_write)(struct Fd *fd, const void *buf, size_t len);
8     int (*dev_close)(struct Fd *fd);
9     int (*dev_stat)(struct Fd *fd, struct Stat *stat);
10    int (*dev_trunc)(struct Fd *fd, off_t length);
11
12    int (*dev_read_map)(struct Fd *fd, void *addr, off_t offset, int perm);
13 };

```

```

1 // lib/file.c
2 static int
3 devfile_read_map(struct Fd *fd, void *addr, off_t offset, int perm)
4 {
5     int r;
6
7     fsipcbuf.read_map.req_fileid = fd->fd_file.id;
8     fsipcbuf.read_map.req_offset = offset;
9     fsipcbuf.read_map.req_perm = perm;
10    if ((r = fsipc(FSREQ_READ_MAP, addr)) < 0)
11        return r;
12
13    return r;
14 }

```

With `read_map()`, we can share the read-only pages during `spawn()` without making a copy for the child.

```

1 // lib/spawn.c map_segment()
2 for (i = 0; i < memsz; i += PGSIZE) {

```

```

3         if (i >= filesz) {
4             // allocate a blank page
5             if ((r = sys_page_alloc(child, (void*) (va + i), perm)) < 0)
6                 return r;
7         } else {
8             // from file
9             if (~perm & PTE_W)
10            {
11                // use read_map() to share read-only pages
12                // between environments
13                if ((r = read_map(fd, UTEMP, fileoffset + i, perm)) < 0)
14                    return r;
15                if ((r = sys_page_map(0, UTEMP, child, (void*) (va + i),
perm)) < 0)
16                    panic("spawn: sys_page_map sharing data: %e", r);
17                sys_page_unmap(0, UTEMP);
18
19                continue;
20            }
21            if ((r = sys_page_alloc(0, UTEMP, PTE_P|PTE_U|PTE_W)) < 0)
22                return r;
23            if ((r = seek(fd, fileoffset + i)) < 0)
24                return r;
25            if ((r = readn(fd, UTEMP, MIN(PGSIZE, filesz-i))) < 0)
26                return r;
27            if ((r = sys_page_map(0, UTEMP, child, (void*) (va + i), perm))
< 0)
28                panic("spawn: sys_page_map data: %e", r);
29            sys_page_unmap(0, UTEMP);
30        }
31    }

```

Sharing library state across fork and spawn

We want share some specific pages e.g. file descriptor with child environment created by `fork()` or `spawn()`. This is implemented by leveraging one of the available bits in PTE to mark pages to share.

Exercise 8

In `lib/fork.c`, we change `duppage()` to map pages with `PTE_SHARE` set directly into the new environment's address space.

```

1 static int
2 duppage(envid_t envid, unsigned pn)
3 {
4     int r;
5
6     // LAB 4: Your code here.
7     pte_t *pte = PGADDR(PDX(UVPT), pn / (PGSIZE >> 2), (pn % (PGSIZE >> 2))
<< 2);
8     void *addr = (void *) (pn * PGSIZE);
9
10    if (((*pte) & (PTE_W | PTE_COW)) && (~(*pte) & PTE_SHARE))

```



```

11     {
12         // writable or copy-on-write page
13         int perm = ((*pte) & PTE_SYSCALL) & (~PTE_W) | PTE_COW;
14         if ((r = sys_page_map(0, addr, envid, addr, perm)))
15         {
16             panic("duppage: %e\n", r);
17         }
18         if ((r = sys_page_map(0, addr, 0, addr, perm)))
19         {
20             panic("duppage: %e\n", r);
21         }
22     }
23     else
24     {
25         // read-only page or shared page
26         int perm = (*pte) & PTE_SYSCALL;
27         if ((r = sys_page_map(0, addr, envid, addr, perm)))
28         {
29             panic("duppage: %e\n", r);
30         }
31     }
32     // panic("duppage not implemented");
33     return 0;
34 }

```

In `lib/spawn.c`, we check every existent page in user space and copy those pages with `PTE_SHARE` set in `copy_shared_pages()`.

```

1 // Copy the mappings for shared pages into the child address space.
2 static int
3 copy_shared_pages(envid_t child)
4 {
5     // LAB 5: Your code here.
6     int r;
7     for (int i = 0; i < PDX(UTOP); i++) if (uvpd[i] & PTE_P)
8     {
9         for (int j = 0; j < NPTENTRIES; j++)
10        {
11            void *addr = PGADDR(i, j, 0);
12            pte_t pte = uvpt[PGNUM(addr)];
13            if ((pte & PTE_P) && (pte & PTE_SHARE))
14            {
15                int perm = pte & PTE_SYSCALL;
16                if ((r = sys_page_map(0, addr, child, addr, perm)) < 0)
17                    return r;
18            }
19        }
20    }
21    return 0;
22 }

```

The keyboard interface

Exercise 9

Now to provide input function for user environments, we dispatch the hardware interrupt from keyboard and serial port to `kbd_intr()` and `serial_intr()` respectively, in which `cons_intr()` is called with corresponding function for requiring data from the hardware. In `cons_intr()`, input characters are stored in a ring buffer which is to be drained by the console file type defined in `lib/console.c`

```
1 // Handle keyboard and serial interrupts.
2 // LAB 5: Your code here.
3 if (tf->tf_trapno == IRQ_OFFSET + IRQ_KBD)
4 {
5     kbd_intr();
6     return;
7 }
8 if (tf->tf_trapno == IRQ_OFFSET + IRQ_SERIAL)
9 {
10     serial_intr();
11     return;
12 }
```

The Shell

Exercise 10

In order to support IO redirection, we need to make the file descriptor ID of the file provided to be 0 or 1. `init` has set up console as file descriptors 0 and 1 and this is inherited by the new environment. So after opening the file we must exploit the `dup(int oldfd, int newfd)` method to close the original descriptor 0 and set it to point to our provided file. Then we close the temporary file descriptor to prevent leak.

```
1 // LAB 5: Your code here.
2 // panic("< redirection not implemented");
3 if ((fd = open(t, O_RDONLY)) < 0) {
4     cprintf("open %s for read: %e", t, fd);
5     exit();
6 }
7 if (fd != 0) {
8     dup(fd, 0);
9     close(fd);
10 }
11 break;
```

Reference

[1] Intel® 64 and IA-32 Architectures Software Developer's Manual, Volume 3A, Section 2.3