

Report for lab5, YiChen Mao

Report for lab5, YiChen Mao

The File System

Exercise 1.
Exercise 2.
Exercise 3.
Exercise 4.
Exercise 5.
Exercise 6.
Exercise 7.
Exercise 8.
Exercise 9.
Exercise 10.
challenge

The File System

主要需要实现的功能，从磁盘中读取块到块缓存中，或从缓存中写回到磁盘中，分配磁盘块，将文件偏移映射到磁盘块中，在 IPC 接口时间读写打开。

Exercise 1.

Task. 修改 env.c 中 env_create 使得赋予环境 I/O 权限。

```
1  在这个 lab 中新加入了一个环境类型 ENV_TYPE_FS，只有环境类型为 ENV_TYPE_FS 才能
   赋予 I/O 权限，ENV_TYPE_USER 则不行。
2  具体修改只需要修改 eflags 对应位，在 mmu.h 中可以找到所有关于 eflags 的宏定义。
3  void
4  env_create(uint8_t *binary, enum EnvType type)
5  {
6      ...
7      if (type == ENV_TYPE_FS){
8          e->env_tf.tf_eflags |= FL_IOPL_MASK;
9      }
10 }
```

Question. 在切换环境的时候是否需要额外操作保证 I/O 权限正确性。

答：我认为是不需要的，因为在切换的过程中，会修改对应的 env 包括其中的 eflags 因此保证了正确。

Exercise 2.

Task. 实现 fs/bc.c 中函数 bc_pgfault 和 flush_block

```
1  磁盘的 page_fault 处理，将数据从磁盘读取到对应的内存
2  static void
3  bc_pgfault(struct UTrapframe *utf)
4  {
5      ...
6      addr = (void*)ROUNDDOWN(addr, BLKSIZE);
7      if ((r = sys_page_alloc(0, addr, PTE_P|PTE_U|PTE_W)) < 0)
8          panic("bc_pgfault: sys_page_alloc (addr) is %x\n", addr);
9      if ((r = ide_read(blockno * BLKSECTS, addr, BLKSECTS)) < 0)
10         panic("bc_pgfault: ide_read (blockno) is %x, (addr) is
11         %x\n", blockno, addr);
12
13     if ((r = sys_page_map(0, addr, 0, addr, uvpt[PGNUM(addr)] &
14     PTE_SYSCALL)) < 0)
15         panic("in bc_pgfault, sys_page_map: %e", r);
16
17     if (bitmap && block_is_free(blockno))
18         panic("reading free block %08x\n", blockno);
19 }
```

```
1  将 cache 中的内容写回到磁盘之中去。
2  条件：当前页被映射了且被修改了，那么需要再写回保证数据正确性。
3  同时修改当前 cache 中的标记位，以防重复写回未修改的数据。
4  void
5  flush_block(void *addr)
6  {
7      uint32_t blockno = ((uint32_t)addr - DISKMAP) / BLKSIZE;
8
9      if (addr < (void*)DISKMAP || addr >= (void*)(DISKMAP + DISKSIZE))
10         panic("flush_block of bad va %08x", addr);
11
12     int r;
```

```

13     addr = (void*)ROUNDDOWN(addr, BLKSIZE);
14     if (va_is_mapped(addr)&&va_is_dirty(addr)){
15         if ((r = ide_write(blockno * BLKSECTS, addr, BLKSECTS))<0)
16             panic("flush_block: ide_write (blockno) is %x, (addr) is
17 %x", blockno, addr);
18         if ((r = sys_page_map(0, addr, 0, addr, uvpt[PGNUM(addr)]&PTE_SYSCALL))
19 <0)
20             panic("flush_block: sys_page_map (addr) is %x\n", addr);
21     }
22     return;
23     panic("flush_block not implemented");
24 }

```

Exercise 3.

Task. 实现 fs/fs 函数 alloc_block()。

```

1  alloc_block 找到一个空闲的磁盘块并分配它。
2  如果未能分配块就 返回 -E_NO_DISK
3
4  从 第 2 个块开始查找空闲块。这里用 block_is_free 判断一个块是否空间，维护了一个
   bitset，第 i 位对应第 i 个块是否为空。
5  映射后修改 bitset 并写回到磁盘中。
6  int
7  alloc_block(void)
8  {
9      // The bitmap consists of one or more blocks. A single bitmap block
10     // contains the in-use bits for BLKBITSIZE blocks. There are
11     // super->s_nblocks blocks in the disk altogether.
12
13     // LAB 5: Your code here.
14     int free_block = 2;
15     while (free_block < super->s_nblocks &&
16 !block_is_free(free_block))free_block++;
17     if (free_block == super->s_nblocks)return -E_NO_DISK;
18     bitmap[free_block/32] &= ~(1<<(free_block%32));
19     flush_block(&bitmap[free_block/32]);
20     return free_block;
21
22     panic("alloc_block not implemented");
23     return -E_NO_DISK;
24 }

```

Exercise 4.

Task. 实现 `file_block_walk` 和 `file_get_block`。

```
1  file_block_walk 在文件 f 中找到第 filebno 个 block 所对应的磁盘块，将对应
   slot 写到 ppdiskbno 中。
2
3  分几类情况讨论：
4  1. 首先 filebno < NDIRECT(10) ，此时直接从当前文件的 f_direct 即可获得对应的指
   针地址。
5  2. 其次 filebno > NDIRECT 此时要到 f_indirect 对应的磁盘块查找地址。如果
   f_indirect 未分配则先分配一个磁盘页。建立映射，并返回新的磁盘页的指针地址。
6  static int
7  file_block_walk(struct File *f, uint32_t filebno, uint32_t **ppdiskbno,
   bool alloc)
8  {
9      // LAB 5: Your code here.
10     int r;
11     if (filebno >= NDIRECT + NINDIRECT)
12         return -E_INVAL;
13     if (filebno >= NDIRECT && f->f_indirect == 0){
14         if (!alloc)return -E_NOT_FOUND;
15         if ((r = alloc_block()) < 0)return r;
16
17         f->f_indirect = r;
18         memset(diskaddr(r),0,BLKSIZE);
19         flush_block(diskaddr(r));
20     }
21     *ppdiskbno = (filebno < NDIRECT)?&f->f_direct[filebno]:
   (uint32_t*)diskaddr(f->f_indirect)+(filebno - NDIRECT);
22     return 0;
23     panic("file_block_walk not implemented");
24 }
```

1 这里找到真正对应的磁盘块。 `file_block_walk` 只是找到文件对应块的磁盘块的地址。
2 相应的，如果该地址未分配对应的磁盘块则分配块，并建立地址和磁盘块之间的映射。
3 最后返回地址对应的磁盘块即可。

```
4  int
5  file_get_block(struct File *f, uint32_t filebno, char **blk)
6  {
```

```

7   int r;
8   uint32_t *pdiskbno;
9   if ((r = file_block_walk(f, filebno, &pdiskbno, true)) < 0)
10      return r;
11   // cprintf("file_get_block: find\n");
12   if (*pdiskbno == 0){
13       if ((r = alloc_block()) < 0)
14           return r;
15       *pdiskbno = r;
16       flush_block(diskaddr(r));
17   }
18   // cprintf("file_get_block: find\n");
19   *blk = diskaddr(*pdiskbno);
20   return 0;
21   panic("file_get_block not implemented");
22 }

```

Exercise 5.

Task. 实现 fs/ servlet.c 中 serve_read

```

1   首先 通过 openfile_lookup 找到对应的 OpenFile。
2   然后调用 file_read 将需要读取的内容放入到 ipc->readRet->ret_buf 中。
3   并修改 offset。
4
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       int r;
16       struct OpenFile *o;
17       if ((r = openfile_lookup(envid, req->req_fileid, &o)) < 0)
18           return r;
19       if ((r = file_read(o->o_file, ret->ret_buf, req->req_n, o->o_fd-
>fd_offset)) < 0)

```

```

20     return r;
21     o->o_fd->fd_offset += r;
22     return r;
23 }
24

```

Exercise 6.

Task. 实现 fs/serv.c 中的 `serve_write` 和 lib/file.c 中的 `devfile_write`。

```

1  和 serve_read 相似, 找到 OpenFile 并写入, 最后修改 offset。
2
3  int
4  serve_write(envid_t envid, struct Fsreq_write *req)
5  {
6      if (debug)
7          cprintf("serve_write %08x %08x %08x\n", envid, req->req_fileid,
8              req->req_n);
9
10     // LAB 5: Your code here.
11     int r;
12     struct OpenFile *o;
13     // cprintf("server_write: %x %x %x %x\n", envid, req->req_fileid, req-
14     >req_buf, &o);
15     if ((r = openfile_lookup(envid, req->req_fileid, &o)) < 0)
16         return r;
17     if ((r = file_write(o->o_file, req->req_buf, (req->req_n>PGSIZE)?
18         PGSIZE: req->req_n, o->o_fd->fd_offset)) < 0)
19         return r;
20     // cprintf("find\n");
21     o->o_fd->fd_offset += r;
22     // cprintf("serve_write fd_offset: %x\n", o->o_fd->fd_offset);
23     return r;
24     panic("serve_write not implemented");
25 }
26

```

```

1  写入至多 n byte 从 buf 到 fd. 可以 仿照 devfile_read.
2  static ssize_t
3  devfile_write(struct Fd *fd, const void *buf, size_t n)

```

```

4 {
5     // Make an FSREQ_WRITE request to the file system server.  Be
6     // careful: fsipcbuf.write.req_buf is only so large, but
7     // remember that write is always allowed to write *fewer*
8     // bytes than requested.
9     // LAB 5: Your code here
10    int r;
11    fsipcbuf.write.req_fileid = fd->fd_file.id;
12    fsipcbuf.write.req_n = (n>sizeof(fsipcbuf.write.req_buf))?
sizeof(fsipcbuf.write.req_buf):n;
13    memcpy(fsipcbuf.write.req_buf,buf,fsipcbuf.write.req_n);
14    // cprintf("write\n");
15    if ((r = fsipc(FSREQ_WRITE,NULL)) < 0)
16        return r;
17    // cprintf("r:%x\n",r);
18    assert(r <= n);
19    assert(r <= PGSIZE);
20    // cprintf("r:%x\n",r);
21    return r;
22    panic("devfile_write not implemented");
23 }

```

Exercise 7.

Task. 实现 sys_env_set_trapframe 函数。

```

1  实现从文件中加载一个子进程，这里只需要实现 设置 tf 即可。
2  由于之前都已经处理好，需要检查用户空间是否合法，然后设置 eflag 即可。
3  static int
4  sys_env_set_trapframe(envid_t envid, struct Trapframe *tf)
5  {
6      int r;
7      struct Env *e;
8      if ((r = envid2env(envid,&e,true)) < 0)
9          return r;
10
11     user_mem_assert(e,(const void*)tf,sizeof(struct Trapframe),PTE_U);
12     tf->tf_eflags = (tf->tf_eflags | FL_IF) & ~FL_IOPL_MASK;
13     tf->tf_cs |= 3;
14     e->env_tf = *tf;
15     return 0;
16     panic("sys_env_set_trapframe not implemented");

```

Exercise 8.

Task. 处理 `duppage()` 中的 `PTE_SHARE` 页，如之前所述，采用直接复制映射的方式。对于 `copy_shared_page`

```

1  在 duppage() 中添加对于 PTE_SHARE 位处理，直接映射，采用 PTE_SYSCALL 作为
    perm。
2  if (pte&PTE_SHARE){
3      if ((r = sys_page_map(envid_parent, (void*)va, env_id,
        (void*)va, PTE_SYSCALL)) < 0)
4          return r;
5      return 0;
6  }
```

```

1  枚举用户空间所在页，对于设置 PTE_SHARE 的页面，同样直接映射，操用 PTE_SYSCALL 作
    为 perm。
2  static int
3  copy_shared_pages(envid_t child)
4  {
5      // LAB 5: Your code here.
6      int r;
7      for (uintptr_t va = 0; va < UTOP; va += PGSIZE){
8          if ((uvpd[PDX(va)] & PTE_P) && (uvpt[PGNUM(va)] & PTE_P) &&
        (uvpt[PGNUM(va)] & PTE_U) && (uvpt[PGNUM(va)] & PTE_SHARE)){
9              if ((r = sys_page_map(0, (void*)va, child, (void*)va, PTE_SYSCALL))
        < 0);
10         }
11     }
12     return 0;
13 }
14
```

Exercise 9.

Task. 调用 `kern/trap.c`，处理 `IRQ_OFFSET+IRQ_KBD` 和 `IRQ_OFFSET+IRQ_SERIAL`。


```

1 kbd_intr() 和 serial_intr() 已经实现好了，因此我们只需要根据 tf_trapno 调用即可。
2 switch (tf->tf_trapno){
3     case IRQ_OFFSET + IRQ_TIMER: {
4         lapic_eoi();
5         sched_yield();
6         return;
7     }
8     case IRQ_OFFSET + IRQ_KBD:{
9         kbd_intr();
10        return;
11    }
12    case IRQ_OFFSET + IRQ_SERIAL: {
13        serial_intr();
14        return;
15    }
16 }

```

Exercise 10.

Task. 实现 shell 对于 < 符号的重定向。

```

1 < 是对于 读文件，因此参数只需要 O_RDONLY。
2 如果 fd 不为 0，则将 fd dup 到 0并关闭原 fd。
3     case '<': // Input redirection
4         // Grab the filename from the argument list
5         if (gettoken(0, &t) != 'w') {
6             cprintf("syntax error: < not followed by word\n");
7             exit();
8         }
9         // Open 't' for reading as file descriptor 0
10        // (which environments use as standard input).
11        // We can't open a file onto a particular descriptor,
12        // so open the file as 'fd',
13        // then check whether 'fd' is 0.
14        // If not, dup 'fd' onto file descriptor 0,
15        // then close the original 'fd'.
16
17        // LAB 5: Your code here.
18        if ((fd = open(t,O_RDONLY))<0){

```

```

19         cprintf("fd open error");
20         exit();
21     }
22     if (fd){
23         dup(fd,0);
24         close(fd);
25     }
26     break;

```

challenge

Challenge! The block cache has no eviction policy. Once a block gets faulted in to it, it never gets removed and will remain in memory forevermore. Add eviction to the buffer cache. Using the PTE_A "accessed" bits in the page tables, which the hardware sets on any access to a page, you can track approximate usage of disk blocks without the need to modify every place in the code that accesses the disk map region. Be careful with dirty blocks.

```

1  #define NBUF 64
2
3  对于 cache 构建一个双向链表结构, unused 表示空闲的 block, used 表示正在使用的
   block。
4  typedef struct MyLink{
5      struct MyLink *p,*s;
6      void *addr;
7  }MyLink;
8  void MyLink_init(MyLink *p){p->p=p->s=p;p->addr=0;}
9  MyLink* MyLink_delete(MyLink *p){p->p->s=p->s;p->s->p=p->p;p->p=p->
   >s=p;return p;}
10 void MyLink_insert(MyLink *p,MyLink *q,void *addr){
11     //cprintf("link %x %x\n",p,q);
12     q->p=p;q->s=p->s;p->s=q->s->p=q;q->addr=addr;
13 }
14 MyLink buf[NBUF],unused,used;
15
16
17 int nbuf=0,t=0;
18
19 将一个 block 移出 cache。

```

```

20 如果该块被修改过，那么要将其写回到磁盘中，并且取消内存映射，并将其从使用链表中删除并
    加入到空闲链表。
21 int buf_remove(MyLink *l){
22     // cprintf("remove %x %x\n",l,l->addr);
23     int r;
24     if (uvpt[PGNUM(l->addr)]&PTE_D)flush_block(l->addr);
25     if (r=sys_page_unmap(0,l->addr),r<0)
26         return r;
27     nbuf--;
28     // cprintf("fffff  %x %x\n",unused.p,l);
29     MyLink_insert(unused.p,MyLink_delete(l),0);
30     return 0;
31 }
32 但 cache 满时，需要将一个 block 移除。
33 选择的策略是 根据是否最近访问过，该标记由 PTE_A 维护，如果一个 block 最近未访问
    过，则将其删除，否则删除一个最早加入 cache 的block。
34 int buf_evict(void){
35     // cprintf("%x %x\n",&used,used.s);
36     int r;
37     for (MyLink *l=used.s;l!=&used;){
38         void *addr=l->addr;
39         if (!(uvpt[PGNUM(addr)]&PTE_A)){
40             //cprintf("%x %x %x %x\n",addr,l,buf,&used);
41             MyLink* tmp = l->s;
42             if((r=buf_remove(l))&&r<0)
43                 return r;
44             l = tmp;
45         }else l = l->s;
46     }
47     if (nbuf==NBUF&&(r=buf_remove(used.s))&&r<0)
48         return r;
49     return 0;
50 }
51 输出函数，能够将所有当前在 cache 中 的块输出。
52 同时输出空间利用率，这里将 PTE_A 标记为 1 的认为是当前使用的 block，而在 cache
    中而未标记 PTE_A 的则认为是未使用 void buf_print_used(void){
53     static int u=0,un=0;
54     for (MyLink *l = used.s;l!=&used;l=l->s){
55         cprintf("%x ",l->addr);
56         if (uvpt[PGNUM(l->addr)]&PTE_A)u++;
57         else un++;
58     }cprintf("\n");
59     cprintf("%d %d %d\n",u,u+un,(int)((float)u/(u+un)*1000000000));

```

```

60 }
61 从 cache 中分配一个 block。首先忽略 前两个 block，因为他们不能被换出。
62 如果 没有空闲 则从 buf 中移除一个。
63 然后从空闲删除一个并加入到使用链表中。存储对应的地址。
64 int buf_alloc(void *addr){
65     // cprintf("alloc %x\n",addr);
66     int r;
67     if (addr<diskaddr(2))return 0;
68     if ((nbuf==NBUF)&&(r=buf_evict())&&r<0)
69         return r;
70     nbuf++;
71     //cprintf("important %x %x\n",&used,&unused);
72     MyLink * l = MyLink_delete(unused.s);
73     MyLink_insert(used.p,l,addr);
74     return 0;
75 }
76 初始化, unused 链表中是所有 cache , 而 used 链表为空。
77 void buf_init(void){
78     MyLink_init(&unused);MyLink_init(&used);
79     for (int i=0;i<NBUF;++i) {
80         MyLink_init(&buf[i]);
81         MyLink_insert(unused.p,&buf[i],0);
82     }
83     /*
84     cprintf("init:\n");
85     for (MyLink *l=unused.s;l!=&unused;l=l->s){
86         cprintf("%x ",l);
87     }
88     cprintf("\n");*/
89 }
90 block 访问, 每次对于 block 的访问都会导致 visit 次数加 1, 当达到一定次数之后, 回
    将 PTE_A 标记清空。
91 int buf_visit(void){
92     int r;
93     if (++t<NBUF)return 0;t=0;
94     int cused = 0,cunused = 0;
95     for (MyLink *l=used.s;l!=&used;l=l->s)cused++;
96     for (MyLink *l=unused.s;l!=&unused;l=l->s)cunused++;
97     // cprintf("%x %x\n",cused,(cused+cunused));
98     for (MyLink *l=used.s;l!=&used;l=l->s){
99         void *addr=l->addr;
100         if (uvpt[PGNUM(addr)]&PTE_A){
101             if (uvpt[PGNUM(addr)]&PTE_D)flush_block(addr);

```

```

102         if
103         ((r=sys_page_map(0,addr,0,addr,uvpt[PGNUM(addr)]&PTE_SYSCALL))&&r<0)
104             return r;
105     }
106     return 0;
107 }
108 这里是对 block_free 的时候需要将对应地址的 block 从 cache 中删除。
109 int buf_delete(void*addr){
110     for (MyLink *l=used.s;l!=&used;l=l->s){
111         if (l->addr==addr){
112             MyLink_insert(unused.p,MyLink_delete(l),0);
113             return 0;
114         }
115     }
116     return -1;
117 }
118

```

```

1  在 bc_init() 中 调用 buf_init() 初始化。
2  void
3  bc_init(void)
4  {
5      struct Super super;
6      set_pgfault_handler(bc_pgfault);
7      check_bc();
8  #ifdef BUF_CACHE_OPEN
9      buf_init();
10 #endif
11     // cache the super block by reading it once
12     memmove(&super, diskaddr(1), sizeof super);
13 }
14
15 在 bc_pgfault 中会将一个块移入内存中, 因此需要在 cache 中加入该 block
16 static void
17 bc_pgfault(struct UTrapframe *utf)
18 {
19     ...
20 #ifdef BUF_CACHE_OPEN
21     if ((r = buf_alloc(addr))&&r < 0)
22         panic("in bc_pgfault, buf_alloc: %e", r);
23 #endif

```

```

24
25 }
26
27 free_block 中需要将其从 cache 中删除。
28 void
29 free_block(uint32_t blockno)
30 {
31     // Blockno zero is the null pointer of block numbers.
32     if (blockno == 0)
33         panic("attempt to free zero block");
34     bitmap[blockno/32] |= 1<<(blockno%32);
35
36 #ifdef BUF_CACHE_OPEN
37     int r;
38     if ((r = sys_page_unmap(0, diskaddr(blockno)))&&r < 0)
39         panic("free_block: %e", r);
40     extern int buf_delete(void*);
41     if ((r=buf_delete(diskaddr(blockno)))&&r<0)
42         panic("free_block: %e", r);
43
44 #endif
45     //cprintf("free_block %x\n",blockno);
46 }

```

测试。用 testfile , testpteshare 和 testshell 文件测试（因为这三个测试的访问 block 较多）

测试数据	NBUF=64	NBUF=32	NBUF=16	NBUF=8	NBUF=4
testfile	37/0/0.3523	69/51/0.4277	69/53/0.5031	70/64/0.6885	214/210/0.9199
testpteshare	11/0/0.8440	11/0/0.8440	11/0/0.8440	13/9/0.8119	18/14/0.8333
testshell	34/0/0.1214	70/44/0.1397	89/74/0.1969	146/139/0.4501	289/285/0.6945

表中第一项数据为cache 中 分配block 次数，第二项为 删除block 的次数，第三项可以认为是空间利用率。

计算方式如下：

用 PTE_A 表示该块是否最近使用过，对于一个随机时刻，记录分配块中最近使用和未使用的块的个数。多次采样取平均值。（这里的采样是对于每次访问就采一次样）

总体而言，利用率随着 BUF 减小而增大，而切换次数也逐渐增多。

(testpteshare 行利用率局部的降低是因为我的程序 NBUF 越少，我刷新 PTE_A 的频率就会越块，因此较小的 BUF 在刷新 PTE_A 之后导致 used 中块大多数变成最近未使用了。)