

lab8 report

lab8 report

[练习0: 填写已有实验](#)

`proc.c:alloc_proc`

`proc.c:do_fork`

[练习1: 完成读文件操作的实现\(需要编码\)](#)

[实现思路](#)

[相关函数](#)

[关键代码](#)

[UNIX的PIPE机制](#)

[练习2: 完成基于文件系统的执行程序机制的实现\(需要编码\)](#)

`load_icode`

[示例执行结果](#)

[UNIX的硬链接和软链接机制](#)

[与参考答案的区别](#)

[练习0](#)

[练习1,2](#)

[知识点](#)

[列出你认为本实验中重要的知识点，以及与对应的OS原理中的知识点，并简要说明你对二者的含义，关系，差异等方面的理解（也可能出现实验中的知识点没有对应的原理知识点）](#)

[列出你认为OS原理中很重要，但在实验中没有对应上的知识点](#)

练习0：填写已有实验

使用 `meld` 填写已有实验。

对之前实验代码的修改：

`proc.c:alloc_proc`

增加对进程相关文件信息的初始化

```
1 | proc->files = NULL;
```

`proc.c:do_fork`

fork时增加对打开文件信息的复制：

```

1  int
2  do_fork(uint32_t clone_flags, uintptr_t stack, struct trapframe *tf) {
3      ...
4      //    3. call copy_mm to dup OR share mm according clone_flag
5      if(copy_mm(clone_flags,proc)!=0){
6          goto bad_fork_cleanup_kstack;
7      }
8
9      // for lab8
10     if(copy_files(clone_flags,proc)!=0){
11         goto bad_fork_cleanup_fs;
12     }
13
14     //    4. call copy_thread to setup tf & context in proc_struct
15     copy_thread(proc,stack,tf);
16     ...
17     fork_out:
18         return ret;
19
20     bad_fork_cleanup_fs: //for LAB8
21         put_files(proc);
22     bad_fork_cleanup_kstack:
23         put_kstack(proc);
24     bad_fork_cleanup_proc:
25         kfree(proc);
26         goto fork_out;
27 }

```

练习1: 完成读文件操作的实现(需要编码)

首先了解打开文件的处理流程，然后参考本实验后续的文件读写操作的过程分析，编写在sfs_inode.c中sfs_io_nolock读文件中数据的实现代码。

请在实验报告中给出设计实现“UNIX的PIPE机制”的概要设方案，鼓励给出详细设计方案

实现思路

在sfs_io_nolock函数中，先计算一些辅助变量，并处理一些特殊情况（比如越界），然后有sfs_buf_op = sfs_rbuf,sfs_block_op = sfs_rblock，设置读取的函数操作。接着进行实际操作，先处理起始的没有对齐到块的部分，再以块为单位循环处理中间的部分，最后处理末尾剩余的部分。每部分中都调用sfs_bmap_load_nolock函数得到blkno对应的inode编号，并调用sfs_rbuf或sfs_rblock函数读取数据（中间部分调用sfs_rblock，起始和末尾部分调用sfs_rbuf），调整相关变量。完成后如果offset + alen > din->fileinfo.size（写文件时会出现这种情况，读文件时不会出现这种情况，alen为实际读写的长度），则调整文件大小为offset + alen并设置dirty变量。

相关函数

`sfs_bmap_load_nolock` 函数将对应 `sfs_inode` 的第 `index` 个索引指向的 `block` 的索引值取出存到相应的指针指向的单元 (`ino_store`)。它调用 `sfs_bmap_get_nolock` 来完成相应的操作。`sfs_rbuf` 和 `sfs_rblock` 函数最终都调用 `sfs_rwblock_nolock` 函数完成操作，而 `sfs_rwblock_nolock` 函数调用 `dop_io->disk0io->disk0read_blks_nolock->ide_read_secs` 完成对磁盘的操作。

关键代码

先处理起始的没有对齐到块的部分

```
1      blkoff = offset % SFS_BLKSIZE;
2      if (blkoff != 0) {
3          size = (nblks != 0) ? (SFS_BLKSIZE - blkoff) : (endpos - offset);
4          if ((ret = sfs_bmap_load_nolock(sfs, sin, blkno, &ino)) != 0) {
5              goto out;
6          }
7          if ((ret = sfs_buf_op(sfs, buf, size, ino, blkoff)) != 0) {
8              goto out;
9          }
10         alen += size;
11         if (nblks == 0) {
12             goto out;
13         }
14         buf += size;
15         blkno ++;
16         nblks --;
17     }
```

再以块为单位循环处理中间的部分

```
1      size = SFS_BLKSIZE;
2      while (nblks != 0) {
3          if ((ret = sfs_bmap_load_nolock(sfs, sin, blkno, &ino)) != 0) {
4              goto out;
5          }
6          if ((ret = sfs_block_op(sfs, buf, ino, 1)) != 0) {
7              goto out;
8          }
9          alen += size;
10         buf += size;
11         blkno ++;
12         nblks --;
13     }
```

最后处理末尾剩余的部分

```

1      if ((size = endpos % SFS_BLKSIZE) != 0) {
2          if ((ret = sfs_bmap_load_nolock(sfs, sin, blkno, &ino)) != 0) {
3              goto out;
4          }
5          if ((ret = sfs_buf_op(sfs, buf, size, ino, 0)) != 0) {
6              goto out;
7          }
8          alen += size;
9      }

```

UNIX的PIPE机制

用先进先出队列做缓冲区。进程写的数据时，若缓冲区没有满，入队，否则将进程阻塞。当缓冲区不是空时，进程可以读数据，出队，否则也将进程阻塞。需要考虑到同步互斥问题，可用信号量或者管程来实现。

练习2: 完成基于文件系统的执行程序机制的实现(需要编码)

改写proc.c中的load_icode函数和其他相关函数，实现基于文件系统的执行程序机制。执行：make qemu。如果能看看到sh用户程序的执行界面，则基本成功了。如果在sh用户界面上可以执行“ls”、“hello”等其他放置在sfs文件系统中的其他执行程序，则可以认为本实验基本成功。

请在实验报告中给出设计实现基于“UNIX的硬链接和软链接机制”的概要设计方案，鼓励给出详细设计方案

load_icode

按照提示将原来从内存读取、拷贝程序的操作转换为从磁盘读取、拷贝。具体步骤见注释。

```

1  static int
2  load_icode(int fd, int argc, char **kargv) {
3      /* LAB8:EXERCISE2 2014011336 HINT:how to load the file with handler fd
4      in to process's memory? how to setup argc/argv?
5
6      * MACROs or Functions:
7      * mm_create          - create a mm
8      * setup_pgdir        - setup pgdir in mm
9      * load_icode_read    - read raw data content of program file
10     * mm_map              - build new vma
11     * pgdir_alloc_page   - allocate new memory for TEXT/DATA/BSS/stack
12     parts
13     * lcr3                - update Page Directory Addr Register -- CR3
14     */
15     /* (1) create a new mm for current process
16     * (2) create a new PDT, and mm->pgdir= kernel virtual addr of PDT
17     * (3) copy TEXT/DATA/BSS parts in binary to memory space of process
18     * (3.1) read raw data content in file and resolve elfhdr

```

```

16      *   (3.2) read raw data content in file and resolve proghdr based on
info in elfhdr
17      *   (3.3) call mm_map to build vma related to TEXT/DATA
18      *   (3.4) callpgdir_alloc_page to allocate page for TEXT/DATA, read
contents in file
19      *           and copy them into the new allocated pages
20      *   (3.5) callpgdir_alloc_page to allocate pages for BSS, memset zero
in these pages
21      * (4) call mm_map to setup user stack, and put parameters into user
stack
22      * (5) setup current process's mm, cr3, reset pgidr (using lcr3 MARCO)
23      * (6) setup uargc and uargv in user stacks
24      * (7) setup trapframe for user environment
25      * (8) if up steps failed, you should cleanup the env.
26      */
27      if (current->mm != NULL) {
28          panic("load_icode: current->mm must be empty.\n");
29      }
30
31      int ret = -E_NO_MEM;
32      struct mm_struct *mm;
33      //(1) create a new mm for current process
34      if ((mm = mm_create()) == NULL) {
35          goto bad_mm;
36      }
37      //(2) create a new PDT, and mm->pgdir= kernel virtual addr of PDT
38      if (setup_pgdir(mm) != 0) {
39          goto bad_pgdir_cleanup_mm;
40      }
41
42      //(3) copy TEXT/DATA/BSS parts in binary to memory space of process
43      struct Page *page;
44
45      //(3.1) read raw data content in file and resolve elfhdr
46      struct elfhdr __elf, *elf = &__elf;
47      if ((ret = load_icode_read(fd, elf, sizeof(struct elfhdr), 0)) != 0) {
48          goto bad_elf_cleanup_pgdir;
49      }
50
51      if (elf->e_magic != ELF_MAGIC) {
52          ret = -E_INVALID ELF;
53          goto bad_elf_cleanup_pgdir;
54      }
55
56      struct proghdr __ph, *ph = &__ph;
57      uint32_t vm_flags, perm, phnum;
58      for (phnum = 0; phnum < elf->e_phnum; phnum++) {
59          //(3.2) read raw data content in file and resolve proghdr based on info
in elfhdr

```

```

60     off_t phoff = elf->e_phoff + sizeof(struct proghdr) * phnum;
61     if ((ret = load_icode_read(fd, ph, sizeof(struct proghdr), phoff))
    != 0) {
62         goto bad_cleanup_mmap;
63     }
64     if (ph->p_type != ELF_PT_LOAD) {
65         continue ;
66     }
67     if (ph->p_filesz > ph->p_memsz) {
68         ret = -E_INVALID_ELF;
69         goto bad_cleanup_mmap;
70     }
71     if (ph->p_filesz == 0) {
72         continue ;
73     }
74     //(3.3) call mm_map to build vma related to TEXT/DATA
75     vm_flags = 0, perm = PTE_U;
76     if (ph->p_flags & ELF_PF_X) vm_flags |= VM_EXEC;
77     if (ph->p_flags & ELF_PF_W) vm_flags |= VM_WRITE;
78     if (ph->p_flags & ELF_PF_R) vm_flags |= VM_READ;
79     if (vm_flags & VM_WRITE) perm |= PTE_W;
80     if ((ret = mm_map(mm, ph->p_va, ph->p_memsz, vm_flags, NULL)) != 0)
    {
81         goto bad_cleanup_mmap;
82     }
83     off_t offset = ph->p_offset;
84     size_t off, size;
85     uintptr_t start = ph->p_va, end, la = ROUNDDOWN(start, PGSIZE);
86
87     ret = -E_NO_MEM;
88
89     end = ph->p_va + ph->p_filesz;
90     //(3.4) callpgdir_alloc_page to allocate page for TEXT/DATA, read
    contents in file
91     //         and copy them into the new allocated pages
92     while (start < end) {
93         if ((page = pgdir_alloc_page(mm->pgdir, la, perm)) == NULL) {
94             ret = -E_NO_MEM;
95             goto bad_cleanup_mmap;
96         }
97         off = start - la, size = PGSIZE - off, la += PGSIZE;
98         if (end < la) {
99             size -= la - end;
100        }
101        if ((ret = load_icode_read(fd, page2kva(page) + off, size,
    offset)) != 0) {
102            goto bad_cleanup_mmap;
103        }
104        start += size, offset += size;

```

```

105     }
106
107     //(3.5) call pgdir_alloc_page to allocate pages for BSS, memset zero in
these pages
108     end = ph->p_va + ph->p_memsz;
109
110     if (start < la) {
111         /* ph->p_memsz == ph->p_filesz */
112         if (start == end) {
113             continue ;
114         }
115         off = start + PGSIZE - la, size = PGSIZE - off;
116         if (end < la) {
117             size -= la - end;
118         }
119         memset(page2kva(page) + off, 0, size);
120         start += size;
121         assert((end < la && start == end) || (end >= la && start ==
la));
122     }
123     while (start < end) {
124         if ((page = pgdir_alloc_page(mm->pgdir, la, perm)) == NULL) {
125             ret = -E_NO_MEM;
126             goto bad_cleanup_mmap;
127         }
128         off = start - la, size = PGSIZE - off, la += PGSIZE;
129         if (end < la) {
130             size -= la - end;
131         }
132         memset(page2kva(page) + off, 0, size);
133         start += size;
134     }
135 }
136 sysfile_close(fd);
137
138 //(4) call mm_map to setup user stack, and put parameters into user
stack
139 vm_flags = VM_READ | VM_WRITE | VM_STACK;
140 if ((ret = mm_map(mm, USTACKTOP - USTACKSIZE, USTACKSIZE, vm_flags,
NULL)) != 0) {
141     goto bad_cleanup_mmap;
142 }
143 assert(pgdir_alloc_page(mm->pgdir, USTACKTOP-PGSIZE , PTE_USER) !=
NULL);
144 assert(pgdir_alloc_page(mm->pgdir, USTACKTOP-2*PGSIZE , PTE_USER) !=
NULL);
145 assert(pgdir_alloc_page(mm->pgdir, USTACKTOP-3*PGSIZE , PTE_USER) !=
NULL);

```

```

146     assert(pgdir_alloc_page(mm->pgdir, USTACKTOP-4*PGSIZE , PTE_USER) !=
NULL);
147
148     //(5) setup current process's mm, cr3, reset pgidr (using lcr3 MARCO)
149     mm_count_inc(mm);
150     current->mm = mm;
151     current->cr3 = PADDR(mm->pgdir);
152     lcr3(PADDR(mm->pgdir));
153
154     //(6) setup uargc and uargv in user stacks
155     uint32_t argv_size=0, i;
156     for (i = 0; i < argc; i++) {
157         argv_size += strlen(kargv[i],EXEC_MAX_ARG_LEN + 1)+1;
158     }
159
160     uintptr_t stacktop = USTACKTOP -
(argv_size/sizeof(long)+1)*sizeof(long);
161     char** uargv=(char **)(stacktop - argc * sizeof(char *));
162
163     argv_size = 0;
164     for (i = 0; i < argc; i++) {
165         uargv[i] = strcpy((char *)(stacktop + argv_size ), kargv[i]);
166         argv_size +=  strlen(kargv[i],EXEC_MAX_ARG_LEN + 1)+1;
167     }
168
169     stacktop = (uintptr_t)uargv - sizeof(int);
170     *(int *)stacktop = argc;
171
172     //(7) setup trapframe for user environment
173     struct trapframe *tf = current->tf;
174     memset(tf, 0, sizeof(struct trapframe));
175     /* LAB5:EXERCISE1 2014011336
176      * should set tf_cs,tf_ds,tf_es,tf_ss,tf_esp,tf_eip,tf_eflags
177      * NOTICE: If we set trapframe correctly, then the user level process
can return to USER MODE from kernel. So
178      *      tf_cs should be USER_CS segment (see memlayout.h)
179      *      tf_ds=tf_es=tf_ss should be USER_DS segment
180      *      tf_esp should be the top addr of user stack (USTACKTOP)
181      *      tf_eip should be the entry point of this binary program
(elf->e_entry)
182      *      tf_eflags should be set to enable computer to produce
Interrupt
183      */
184     tf->tf_cs = USER_CS;
185     tf->tf_ds = tf->tf_es = tf->tf_ss = USER_DS;
186     tf->tf_esp = stacktop;
187     tf->tf_eip = elf->e_entry;
188     tf->tf_eflags = FL_IF;
189     ret = 0;

```



```

190
191     //(8) if up steps failed, you should cleanup the env.
192 out:
193     return ret;
194 bad_cleanup_mmap:
195     exit_mmap(mm);
196 bad_elf_cleanup_pgdir:
197     put_pgdir(mm);
198 bad_pgdir_cleanup_mm:
199     mm_destroy(mm);
200 bad_mm:
201     goto out;
202 }

```

示例执行结果

ls

```

1  ls
2  @ is  [directory] 2(hlinks) 23(blocks) 5888(bytes) : @'.'
3  [d]   2(h)        23(b)    5888(s)   .
4  [d]   2(h)        23(b)    5888(s)   ..
5  [-]   1(h)        10(b)    40224(s)  testbss
6  [-]   1(h)        10(b)    40204(s)  badsegment
7  [-]   1(h)        10(b)    40204(s)  faultread
8  [-]   1(h)        10(b)    40360(s)  ls
9  [-]   1(h)        10(b)    40200(s)  softint
10 [-]   1(h)        10(b)    40292(s)  priority
11 [-]   1(h)        10(b)    40332(s)  waitkill
12 [-]   1(h)        10(b)    40220(s)  divzero
13 [-]   1(h)        10(b)    40192(s)  pgdir
14 [-]   1(h)        10(b)    40204(s)  sleepkill
15 [-]   1(h)        10(b)    40252(s)  forktree
16 [-]   1(h)        10(b)    40208(s)  faultreadkernel
17 [-]   1(h)        10(b)    40200(s)  badarg
18 [-]   1(h)        10(b)    40200(s)  hello
19 [-]   1(h)        10(b)    40220(s)  sleep
20 [-]   1(h)        10(b)    40228(s)  forktest
21 [-]   1(h)        10(b)    40224(s)  exit
22 [-]   1(h)        10(b)    40304(s)  matrix
23 [-]   1(h)        10(b)    40196(s)  spin
24 [-]   1(h)        10(b)    40200(s)  yield
25 [-]   1(h)        11(b)    44508(s)  sh
26 lsdir: step 4

```

forktest

```
1 $ forktest
2 main-loop: WARNING: I/O thread spun for 1000 iterations
3 I am child 31
4 I am child 30
5 I am child 29
6 I am child 28
7 I am child 27
8 I am child 26
9 I am child 25
10 I am child 24
11 I am child 23
12 I am child 22
13 I am child 21
14 I am child 20
15 I am child 19
16 I am child 18
17 I am child 17
18 I am child 16
19 I am child 15
20 I am child 14
21 I am child 13
22 I am child 12
23 I am child 11
24 I am child 10
25 I am child 9
26 I am child 8
27 I am child 7
28 I am child 6
29 I am child 5
30 I am child 4
31 I am child 3
32 I am child 2
33 I am child 1
34 I am child 0
35 forktest pass.
```

UNIX的硬链接和软链接机制

- 硬链接：在文件描述符中加入引用标记和引用文件指针，标记位为1表示硬链接，指针指向链接的文件。删除文件时进行判断，如果被引用文件的引用次数标记为0则删除该文件。
- 软链接：直接拷贝文件对应的 `inode` 信息。

与参考答案的区别

练习0

参考答案在 `copy_mm` 失败时会跳到 `bad_fork_cleanup_fs`，而在 `copy_files` 失败的时候会跳到 `bad_fork_cleanup_kstack`，明显写反了，希望能改过来。

练习1,2

基本相同

知识点

列出你认为本实验中重要的知识点，以及与对应的OS原理中的知识点，并简要说明你对二者的含义，关系，差异等方面的理解（也可能出现实验中的知识点没有对应的原理知识点）

- 文件系统架构。原理中很清楚的文件系统架构到了具体ucore的实现中则需要好好捋一捋。我自下而上的分析了文件系统，收获很大。
- 进程的文件管理。打开文件，读写操作在实验中都有体现。

列出你认为OS原理中很重要，但在实验中没有对应上的知识点

- 文件系统的挂载
- 路径访问