

# lab8实验报告

组长：谢雯菲 学号：2110803

组员：吴静 学号：2113285

组员：李浩天 学号：2110133

## 练习0：填写已有实验

本实验依赖实验1/2/3/4/5/6/7。请把你做的实验1/2/3/4/5/6/7的代码填入本实验中代码中有“LAB1”/“LAB2”/“LAB3”/“LAB4”/“LAB5”/“LAB6” /“LAB7”的注释相应部分。并确保编译通过。注意：为了能够正确执行lab8的测试应用程序，可能需对已完成的实验1/2/3/4/5/6/7的代码进行进一步改进。

### alloc\_proc() 函数

在 alloc\_proc() 函数中需要对文件结构的指针进行初始化，即将 proc->filesp 赋值为 NULL。

修改后的 alloc\_proc() 函数如下所示：

```
1 static struct proc_struct *
2 alloc_proc(void) {
3     struct proc_struct *proc = kmalloc(sizeof(struct proc_struct));
4     if (proc != NULL) {
5         //LAB4:EXERCISE1 YOUR CODE
6         proc->state = PROC_UNINIT; //设置进程为未初始化状态
7         proc->pid = -1;             //未初始化的进程id=-1
8         proc->runs = 0;             //初始化时间片
9         proc->kstack = 0;           //初始化内存栈的地址
10        proc->need_resched = 0;     //是否需要调度设为不需要
11        proc->parent = NULL;        //置空父节点
12        proc->mm = NULL;            //置空虚拟内存
13        memset(&(proc->context), 0, sizeof(struct context)); //初始化上下文
14        proc->tf = NULL;            //中断帧指针设置为空
15        proc->cr3 = boot_cr3;       //页目录设为内核页目录表的基址
16        proc->flags = 0;            //初始化标志位
17        memset(proc->name, 0, PROC_NAME_LEN); //置空进程名
18        //LAB5 YOUR CODE : (update LAB4 steps)
19        proc->wait_state = 0;       //初始化进程等待状态
20        proc->cptr = proc->optr = proc->yptr = NULL; //进程相关指针初始化
21        //LAB6 YOUR CODE : (update LAB5 steps)
22        proc->rq = NULL;
23        proc->run_link.prev = proc->run_link.next = NULL;
24        proc->time_slice = 0;
25        proc->lab6_run_pool.left = proc->lab6_run_pool.right = proc->
26        >lab6_run_pool.parent = NULL;
27        proc->lab6_stride = 0;
28        proc->lab6_priority = 0;
29        //LAB8 YOUR CODE : (update LAB6 steps)
30        proc->filesp = NULL;
31    }
32    return proc;
33 }
```

## proc\_run() 函数

在 `proc_run()` 函数中，当进程进行上下文切换，即切换到另一个进程运行时，由于不同的进程可能有不同的页表，其虚拟地址到物理地址的映射关系可能会发生变化。因此，在进行进程切换时，需要刷新TLB，以确保TLB中的映射关系与当前运行的进程的页表保持一致。

根据提示，修改后的 `proc_run()` 函数如下所示：

```

1  void
2  proc_run(struct proc_struct *proc) {
3      if (proc != current) {
4          // LAB4:EXERCISE3 YOUR CODE
5          bool intr_flag;
6          struct proc_struct *prev = current;
7          //用于标识当前进程的进程控制块
8          struct proc_struct *next = proc;
9          //用于标识要切换的进程的进程控制块
10         local_intr_save(intr_flag);
11         //确保在调度函数执行期间，不会被中断打断
12         {
13             current = proc;
14             //将当前运行的进程设置为要切换过去的进程
15             lcr3(next->cr3);
16             //将页表换成新进程的页表
17             switch_to(&(prev->context), &(next->context));
18             //使用switch_to切换到新进程
19         }
20         local_intr_restore(intr_flag);
21         //恢复中断
22         //LAB8 YOUR CODE : (update LAB4 steps)
23         flush_tlb();
24     }
25 }
```

## 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，填写在 `kern/fs/sfs/sfs_inode.c` 中的 `sfs_io_nolock()` 函数，实现读文件中数据的代码。

## 打开文件的处理流程

通过学习得知ucore的文件系统总体如下图所示：



一个用户态程序想要执行打开文件的操作的流程为：用户态程序->系统库函数->内核态系统调用->文件系统抽象层处理->SFS文件系统层处理。

## 用户态程序调用系统库函数

在 `user\libs\file.c` 代码中定义了用户态程序对文件可以调用的系统库函数。当用户态程序想要打开文件时，会调用 `open()` 函数，该函数进一步调用 `sys_open()` 的系统库函数，从而进一步调用 `syscall` 引起系统调用进入到内核态。

## 内核态处理

内核调用 `sys_open()` 函数，进一步调用 `sysfile_open()`，完成把位于用户空间的文件路径字符串拷贝到内核空间的操作，最后调用了 `file_open()` 函数进入文件系统抽象层的处理。

## 文件系统抽象层处理

`file_open()` 函数给即将打开的文件分配一个 `file` 数据结构的变量，这个变量其实是 `current->fs_struct->filemap[]` 数组中的一个还没用于打开文件的空闲元素。这个元素的索引值最后会返回到用户进程并赋值给变量 `fd`。也就是说这个函数给当前用户进程分配了一个 `file` 数据结构的变量，还没有找到对应的文件索引节点。

随后，`file_open()` 函数调用 `vfs_open()` 函数找到 `path` 指出的文件所对应的基于 `inode` 数据结构的在抽象文件系统中的索引节点。该函数完成了以下两件事情：

1. 通过 `vfs_lookup()` 找到 `path` 对应文件的 `inode`
2. 调用 `vop_open()` 函数打开文件

其中，`vfs_lookup()` 函数是一个针对目录的操作函数，它会调用 `vop_lookup()` 函数来找到SFS文件系统下的目录下的文件。

在这个流程中，有大量以vop开头的函数，它们都通过一些宏和函数的转发，最后变成对inode结构体里的inode\_ops结构体的“成员函数”的调用，也就是会进一步调用对sfs文件系统操作的函数。

## SFS文件系统层的处理流程

在 `kern\fs\sfs\sfs_inode.c` 中的 `sfs_node_dirops` 变量定义了 `.vop_lookup = sfs_lookup`，所以到了SFS文件系统层时会调用 `sfs_lookup()` 函数。

`sfs_lookup()` 函数会解析传入的目录DIR的路径，以“/”为分隔符，从左至右逐一分解path获得各个子目录和最终文件对应的inode节点。其调用 `sfs_lookup_once()` 函数查找根目录下的文件sfs\_filetest1对应的inode节点。当无法分解path后，就意味着找到了对应的inode节点。

`sfs_lookup_once()` 函数进一步调用 `sfs_dirent_search_nolock()` 函数来查找与路径名匹配的目录项，如果找到目录项，则根据目录项中记录的inode所处的数据块索引值找到路径名对应的SFS磁盘inode，读入SFS磁盘inode对的内容，创建SFS内存inode（将索引值存到slot，将磁盘inode号存到ino\_store）。

之后，文件相应的索引节点值会传回文件系统抽象层中的 `file_open()` 函数，该函数之后会完善file结构的变量，并将要打开的文件对应的文件结构的变量状态设为“打开”。

## 读写文件的处理流程

读写文件总体的流程与打开文件类似：用户态程序->系统库函数->内核态系统调用->文件系统抽象层处理->SFS文件系统层处理。

### 用户态程序调用系统库函数

拿读文件举例，用户态程序调用系统库函数 `read()`，该函数进一步调用 `sys_read()` 函数，从而调用 `syscall` 引起系统调用进入到内核态。

### 内核态处理

到了内核态以后，通过中断处理例程，会调用到 `sys_read()` 内核函数，并进一步调用 `sysfile_read()` 内核函数，进入到文件系统抽象层处理流程完成进一步读文件的操作。

### 文件系统抽象层处理

`sysfile_read()` 函数循环读取文件，每次读取buffer大小。每次循环中先检查剩余部分大小，若其小于4096字节，则只读取剩余部分的大小。然后调用 `file_read()` 函数将文件内容读取到buffer中，alen为实际大小。

`file_read()` 函数首先调用 `fd2file()` 函数找到对应的file结构，并检查是否可读。调用 `filemap_acquire()` 函数使打开这个文件的计数加1。调用 `vop_read()` 函数将文件内容读到iob中。调整文件指针偏移量pos的值，使其向后移动实际读到的字节数 `iobuf_used(iob)`。最后调用 `filemap_release()` 函数使打开这个文件的计数减1，若打开计数为0，则释放file。

## SFS文件系统层的处理流程

在上述文件系统抽象层中调用的 `vop_read()` 函数实际上是对 `sfs_read()` 函数的包装。在 `sfs_read()` 函数中，调用了 `sys_io()` 函数。`sfs_io()` 函数先找到inode对应的sfs和sin，然后调用 `sfs_io_nolock()` 函数进行读取文件操作，最后调用 `iobuf_skip()` 函数调整iobuf的指针。

最后函数返回后，文件系统抽象层的 `sysfile_read()` 函数继续调用 `copy_to_user()` 函数将读到的内容拷贝到用户的内存空间中，调整各变量以进行下一次循环读取，直至指定长度读取完成。最后函数调用层层返回至用户程序，用户程序收到了读到的文件内容。

## sfs\_io\_nolock() 函数实现

在 `sfs_io_nolock()` 函数中需要完成以下操作：

1. 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。
2. 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。
3. 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。
4. 完成后如果 `offset + alen > din->fileinfo.size`（即导致文件大小增加，只有写文件时会出现这种情况），则调整文件大小为 `offset + alen` 并设置dirty变量。

根据提示，需要用到 `sfs_bmap_load_nolock()`、`sfs_buf_op()` 和 `sfs_block_op()` 函数。

## sfs\_bmap\_load\_nolock() 函数

该函数的源码及注释如下所示：

```
1  /* sfs_bmap_load_nolock - 根据目录的 inode 和 inode 中块的逻辑索引，查找磁盘块的编号。
2  * @sfs:      SFS 文件系统
3  * @sin:      内存中的 SFS inode
4  * @index:    inode 中磁盘块的逻辑索引
5  * @ino_store: 存储磁盘块编号的指针
6  */
7  static int
8  sfs_bmap_load_nolock(struct sfs_fs *sfs, struct sfs_inode *sin, uint32_t index,
9                      uint32_t *ino_store)
10 {
11     // 获取指向该 inode 的磁盘表示的指针
12     struct sfs_disk_inode *din = sin->din;
13
14     // 断言所请求的索引不超过当前 inode 的块数
15     assert(index <= din->blocks);
16
17     int ret;
18     uint32_t ino;
19     // 确定是否需要创建新块（写入新的块？）
20     bool create = (index == din->blocks);
```

```

20
21 // 获取给定索引处的块映射，如果需要创建块，则创建新块
22 if ((ret = sfs_bmap_get_nolock(sfs, sin, index, create, &ino)) != 0)
23 {
24     return ret;
25 }
26
27 // 断言获取到的块号确实在使用中
28 assert(sfs_block_inuse(sfs, ino));
29
30 // 如果需要创建块，则增加当前 inode 的块数
31 if (create)
32 {
33     din->blocks++;
34 }
35
36 // 如果 ino_store 指针不为空，则将获取到的块号存储到指针指向的位置
37 if (ino_store != NULL)
38 {
39     *ino_store = ino;
40 }
41
42 // 返回成功加载块映射的结果
43 return 0;
44 }

```

该函数会将对应sfs\_inode的第index个索引指向的block块的索引值取出存到相应的指针指向的单元(ino\_store)。它进一步调用 sfs\_bmap\_get\_nolock() 函数来完成相应的操作。

## sfs\_buf\_op() 和 sfs\_block\_op() 函数

这两个函数实际上是函数的指针，会在 sfs\_bmap\_load\_nolock() 函数的前半部分被赋值。如果是读操作，sfs\_buf\_op() 函数会被赋值为 sfs\_rbuf() 函数，sfs\_block\_op() 函数会被赋值为 sfs\_rblock() 函数；如果是写操作，sfs\_buf\_op() 函数会被赋值为 sfs\_wbuf() 函数，sfs\_block\_op() 函数会被赋值为 sfs\_wblock() 函数。

### sfs\_rbuf() 函数和 sfs\_rblock() 函数

拿读操作举例，sfs\_rbuf() 函数的代码及注释如下所示：

```

1  /*
2   * sfs_rbuf - 用于读取一个磁盘块的基本块级I/O例程（非块对齐和非块I/O），使用 sfs-
   >sfs_buffer,
3   *           并在Rd/Wr磁盘块时进行互斥处理的锁保护
4   * @sfs:     将要进行处理的 sfs_fs 结构体
5   * @buf:     用于读取的缓冲区
6   * @len:     需要读取的长度
7   * @blkno:   磁盘块的编号
8   * @offset:  磁盘块内容中的偏移量
9   */
10 int

```

```

11 sfs_rbuf(struct sfs_fs *sfs, void *buf, size_t len, uint32_t blkno, off_t
offset) {
12     //使用断言确保参数合法性，即偏移量在有效范围内
13     assert(offset >= 0 && offset < SFS_BLKSIZE && offset + len <= SFS_BLKSIZE);
14     int ret;
15     //锁定SFS文件系统，确保在执行I/O操作时互斥
16     lock_sfs_io(sfs);
17     {
18         //调用 sfs_rwblock_nolock 函数来读取指定磁盘块的内容到 SFS 文件系统结构体中的缓冲
区 sfs_buffer。
19         //如果读取成功（返回值为0），则将缓冲区的内容从指定偏移量开始复制到目标缓冲区 buf
中。
20         if ((ret = sfs_rwblock_nolock(sfs, sfs->sfs_buffer, blkno, 0, 1)) == 0)
        {
21             memcpy(buf, sfs->sfs_buffer + offset, len);
22         }
23     }
24     //解锁
25     unlock_sfs_io(sfs);
26     return ret;
27 }

```

sfs\_rblock() 函数进一步调用 sfs\_rwblock() 函数，该函数代码及注释如下所示：

```

1  /*
2   * sfs_rwblock - 用于读取/写入 N 个磁盘块的基本块级 I/O 例程，
3   *               并在 Rd/Wr 磁盘块时进行互斥处理的锁保护
4   * @sfs:        将要进行处理的 sfs_fs 结构体
5   * @buf:         用于读取/写入的缓冲区
6   * @blkno:       磁盘块的编号
7   * @nblks:       读取/写入的磁盘块数量
8   * @write:       布尔值：读取 - 0 或写入 - 1
9   */
10 static int
11 sfs_rwblock(struct sfs_fs *sfs, void *buf, uint32_t blkno, uint32_t nblks, bool
write) {
12     int ret = 0;
13     //锁定SFS文件系统
14     lock_sfs_io(sfs);
15     {
16         //循环处理每个磁盘块，直到读取/写入完指定数量的磁盘块
17         while (nblks != 0) {
18             //调用 sfs_rwblock_nolock 函数来读取或写入一个磁盘块的内容。
19             //如果读取/写入操作失败（返回值不为0），则退出循环。
20             if ((ret = sfs_rwblock_nolock(sfs, buf, blkno, write, 1)) != 0) {
21                 break;
22             }
23             //更新磁盘块编号、剩余磁盘块数量和缓冲区指针，以处理下一个磁盘块。
24             blkno ++, nblks --;
25             buf += SFS_BLKSIZE;
26         }

```



```

27     }
28     //解锁
29     unlock_sfs_io(sfs);
30     return ret;
31 }

```

这两个函数最终都调用 `sfs_rwblock_nolock()` 函数完成操作，而 `sfs_rwblock_nolock()` 函数调用 `dop_io -> disk0_io -> disk0_read_blks_nolock -> ide_read_secs` 完成对磁盘的操作。

## sfs\_io\_nolock() 函数代码实现

根据提示，函数主要分为三个部分实现：

1. 如果偏移量与第一个块不对齐，从偏移位置读取/写入一些内容直到第一个块的末尾
2. 读/写对齐的块
3. 如果结束位置与最后一个块不对齐，从最后一个块的开始位置读取/写入一些内容到最后一个块的 `(endpos % SFS_BLKSIZE)` 处

函数代码最终的实现及注释如下所示：

```

1  static int
2  sfs_io_nolock(struct sfs_fs *sfs, struct sfs_inode *sin, void *buf, off_t
   offset, size_t *alenp, bool write)
3  {
4      .....
5
6      // LAB8:EXERCISE1 YOUR CODE HINT: 2110803
7
8      //(1)blkoff为非对齐的起始块中需要操作的偏移量
9      //判断是否需要操作
10     if ((blkoff = offset % SFS_BLKSIZE) != 0){
11         //得到起始块中要进行操作的数据长度
12         //nbks为0说明为最后一块
13         //如果不是最后一块计算：块大小-操作偏移量
14         //如果是最后一块计算：结束长度-总偏移量
15         size = (nbks != 0) ? (SFS_BLKSIZE - blkoff) : (endpos - offset);
16         //获取这些数据块对应到磁盘上的数据块的inode号
17         if((ret = sfs_bmap_load_nolock(sfs, sin, blkno, &ino)) != 0){
18             goto out;
19         }
20         //对缓冲区进行读或写操作
21         if((ret = sfs_buf_op(sfs, buf, size, ino, blkoff)) != 0){
22             goto out;
23         }
24         //已经完成读写的数据长度
25         alen += size;
26         //如果这是最后操作的块，结束
27         if(nbks == 0){
28             goto out;
29         }
30         //否则，更新缓冲区

```



```

31     buf += size;
32     blkno++;
33     nblks--;
34 }
35 // (2) 读取起始块后对齐块的数据
36 //将中间部分的数据分为一块一块的大小，一块一块操作
37 size = SFS_BLKSIZE;
38 //从起始块后一块开始到对齐的最后一块
39 while(nblks != 0){
40     //获取磁盘上块的编号
41     if((ret = sfs_bmap_load_nolock(sfs, sin, blkno, &ino)) != 0){
42         goto out;
43     }
44     //对数据块进行读或写操作
45     if((ret = sfs_block_op(sfs, buf, ino, 1)) != 0){
46         goto out;
47     }
48     //更新
49     alen += size;
50     buf += size;
51     blkno ++;
52     nblks --;
53 }
54 // (3) 如果有非对齐操作的最后一块
55 if((size = endpos % SFS_BLKSIZE) != 0){
56     //获取磁盘上的块号
57     if((ret = sfs_bmap_load_nolock(sfs, sin, blkno, &ino)) != 0){
58         goto out;
59     }
60     //缓冲区读写操作
61     if((ret = sfs_buf_op(sfs, buf, size, ino, 0)) != 0){
62         goto out;
63     }
64     //更新
65     alen += size;
66 }
67
68 out:
69     .....
70
71     return ret;
72 }

```

## 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）

本次练习将从提示和lab5 `load_icode` 代码入手，提示如下：

```

1     /* (1) create a new mm for current process
2     * (2) create a new PDT, and mm->pgdir= kernel virtual addr of PDT
3     * (3) copy TEXT/DATA/BSS parts in binary to memory space of process
4     *     (3.1) read raw data content in file and resolve elfhdr

```

```

5      *      (3.2) read raw data content in file and resolve proghdr based on info
      in elfhdr
6      *      (3.3) call mm_map to build vma related to TEXT/DATA
7      *      (3.4) call pgdir_alloc_page to allocate page for TEXT/DATA, read
      contents in file
8      *      and copy them into the new allocated pages
9      *      (3.5) call pgdir_alloc_page to allocate pages for BSS, memset zero in
      these pages
10     *      (4) call mm_map to setup user stack, and put parameters into user stack
11     *      (5) setup current process's mm, cr3, reset pgidr (using lcr3 MARCO)
12     *      (6) setup uargc and uargv in user stacks
13     *      (7) setup trapframe for user environment
14     *      (8) if up steps failed, you should cleanup the env.
15     */

```

## (1) create a new mm for current process

调用 `mm_create` 函数来申请进程的内存管理数据结构mm所需内存空间，并对mm进行初始化。

代码如下：

```

1      struct mm_struct *mm;
2      if ((mm = mm_create()) == NULL)
3      {
4          goto bad_mm;
5      }

```

首先使用 `mm_create` 函数创建了一个新的 `mm_struct` 结构体，同时进行了错误检查，如果最后创建的结构体返回值为 `null`，则报错。

`mm_create` 函数如下：

```

1  // mm_create - alloc a mm_struct & initialize it.
2  struct mm_struct *
3  mm_create(void) {
4      struct mm_struct *mm = kmalloc(sizeof(struct mm_struct));
5
6      if (mm != NULL) {
7          list_init(&(mm->mmap_list));
8          mm->mmap_cache = NULL;
9          mm->pgdir = NULL;
10         mm->map_count = 0;
11
12         if (swap_init_ok) swap_init_mm(mm);
13         else mm->sm_priv = NULL;
14
15         set_mm_count(mm, 0);
16         sem_init(&(mm->mm_sem), 1);
17     }
18     return mm;
19 }

```

该函数具体作用是分配一个 `mm_struct` 结构体并且进行初始化。

## (2) create a new PDT, and `mm->pgdir`= kernel virtual addr of PDT

调用 `setup_pgdir` 来申请一个页目录表所需的一个页大小的内存空间，并把描述 `ucore` 内核虚空间映射的内核页表的内容拷贝到此新目录表中，最后让 `mm->pgdir` 指向此页目录表。

代码如下：

```
1     if (setup_pgdir(mm) != 0)
2     {
3         goto bad_pgdir_cleanup_mm;
4     }
```

调用 `setup_pgdir` 函数创建一个新的页目录表，并将 `mm->pgdir` 设置为页目录表的内核虚拟地址。

`setup_pgdir` 函数如下：

```
1  // setup_pgdir - alloc one page as PDT
2  static int
3  setup_pgdir(struct mm_struct *mm)
4  {
5      struct Page *page;
6      if ((page = alloc_page()) == NULL)
7      {
8          return -E_NO_MEM;
9      }
10     pde_t *pgdir = page2kva(page);
11     memcpy(pgdir, boot_pgdir, PGSIZE);
12
13     mm->pgdir = pgdir;
14     return 0;
15 }
```

该函数的作用是分配一个页面作为页目录表。

## (3) copy TEXT/DATA/BSS parts in binary to memory space of process

将磁盘中的文件加载到内存中。

这一步是要将将二进制文件中的 `TEXT/DATA/BSS` 部分复制到进程的内存空间；

注意这里和lab5有一点不同，lab5的 `load_icode` 中传递了二进制文件 `binary` 和大小 `size`，但是本次实验传递的是文件描述符 `fd`，参数数量 `argc` 和参数列表 `kargv`，所以在这一步的具体实现需要做变动。

### (3.0) preparations

在做具体实现前，首先对 `elf` 文件头和程序段头做基础了解。

#### elf文件头格式

```
1  #define ELF_MAGIC    0x464C457FU           // "\x7FELF" in little endian
2
3  /* file header */
4  struct elfhdr {
5      uint32_t e_magic;      // must equal ELF_MAGIC
6      uint8_t e_elf[12];
7      uint16_t e_type;       // 1=relocatable, 2=executable, 3=shared object,
                             // 4=core image
8      uint16_t e_machine;    // 3=x86, 4=68K, etc.
9      uint32_t e_version;    // file version, always 1
10     uint64_t e_entry;       // entry point if executable
11     uint64_t e_phoff;       // file position of program header or 0
12     uint64_t e_shoff;       // file position of section header or 0
13     uint32_t e_flags;       // architecture-specific flags, usually 0
14     uint16_t e_ehsize;      // size of this elf header
15     uint16_t e_phentsize;   // size of an entry in program header
16     uint16_t e_phnum;       // number of entries in program header or 0
17     uint16_t e_shentsize;   // size of an entry in section header
18     uint16_t e_shnum;       // number of entries in section header or 0
19     uint16_t e_shstrndx;    // section number that contains section name strings
20 };
```

下面是对其各个字段的解释：

- `e_magic`: ELF 文件标识，必须等于 `ELF_MAGIC`。
- `e_elf[12]`: ELF 文件类型的额外标识信息。
- `e_type`: 指定 ELF 文件类型
  - 1 表示可重定位文件
  - 2 表示可执行文件
  - 3 表示共享对象
  - 4 表示核心转储文件。
- `e_machine`: 指定了目标体系结构的类型
  - 3 代表 x86
  - 4 代表 68K
- `e_version`: 文件版本号，通常为 1
- `e_entry`: 如果是可执行文件，则表示程序的入口点
- `e_phoff`: 程序头表的文件偏移量，如果没有程序头表则为 0
- `e_shoff`: 节头表的文件偏移量，如果没有节头表则为 0
- `e_flags`: 体系结构相关的标志，通常为 0

- `e_ehsize`: ELF 头的大小 (字节)
- `e_phentsize`: 程序头表中每个条目的大小
- `e_phnum`: 程序头表中的条目数, 如果没有程序头表则为 0
- `e_shentsize`: 节头表中每个条目的大小
- `e_shnum`: 节头表中的条目数, 如果没有节头表则为 0
- `e_shstrndx`: 包含节名称字符串的节的索引号

## 程序段头文件格式

```

1  /* program section header */
2  struct proghdr {
3      uint32_t p_type;    // loadable code or data, dynamic linking info,etc.
4      uint32_t p_flags;  // read/write/execute bits
5      uint64_t p_offset; // file offset of segment
6      uint64_t p_va;     // virtual address to map segment
7      uint64_t p_pa;     // physical address, not used
8      uint64_t p_filesz; // size of segment in file
9      uint64_t p_memsz;  // size of segment in memory (bigger if contains bss)
10     uint64_t p_align;  // required alignment, invariably hardware page size
11 };
12
13 /* values for Proghdr::p_type */
14 #define ELF_PT_LOAD                1
15
16 /* flag bits for Proghdr::p_flags */
17 #define ELF_PF_X                    1
18 #define ELF_PF_W                    2
19 #define ELF_PF_R                    4

```

以下是对每个字段的解释:

- `p_type`: 指定段的类型, 可能是可加载的代码或数据, 动态链接信息等
  - 程序中给了一种类型 `ELF_PT_LOAD`, 表示可加载
- `p_flags`: 标识段的属性
  - `ELF_PF_R` 表示读
  - `ELF_PF_W` 表示写
  - `ELF_PF_X` 表示可执行
- `p_offset`: 段在文件中的偏移量, 指示段在 ELF 文件中的位置
- `p_va`: 段在虚拟地址空间中的起始地址, 用于映射到进程的内存
- `p_pa`: 物理地址, 通常不使用
- `p_filesz`: 段在文件中的大小, 表示段在 ELF 文件中占据的空间大小
- `p_memsz`: 段在内存中的大小, 表示段加载到内存后占据的空间大小, 如果包含未初始化数据 (BSS), 则通常会比文件大小要大
- `p_align`: 段的对齐要求, 通常等于硬件页面大小

接下来，由于本次步骤中的几个步骤都会用上 `load_icode_read` 函数，首先先来解析这个函数的作用。

## `load_icode_read()`

函数原型如下：

```
1 static int
2 load_icode_read(int fd, void *buf, size_t len, off_t offset)
3 {
4     int ret;
5     if ((ret = sysfile_seek(fd, offset, LSEEK_SET)) != 0)
6     {
7         return ret;
8     }
9     if ((ret = sysfile_read(fd, buf, len)) != len)
10    {
11        return (ret < 0) ? ret : -1;
12    }
13    return 0;
14 }
```

该函数接收四个参数：

- 文件描述符 `fd`
- 指向缓冲区的指针 `buf`
- 读取的长度 `len`
- 读取的偏移量 `offset`。

该函数的作用是从文件中读取数据并将其存储到内存缓冲区`buff`中，具体步骤如下：

1. 使用 `sysfile_seek` 函数将文件描述符 `fd` 设置到指定的 `offset`

`sysfile` 函数：

```
1 //kern/fs/sysfile.c
2 int sysfile_seek(int fd, off_t pos, int whence)
3 {
4     return file_seek(fd, pos, whence);
5 }
```

`file_seek` 函数如下：

```
1 // seek file
2 int file_seek(int fd, off_t pos, int whence)
3 {
4     struct stat __stat, *stat = &__stat;
5     int ret;
6     struct file *file;
7     if ((ret = fd2file(fd, &file)) != 0)
8     {
9         return ret;
```

```

10     }
11     fd_array_acquire(file);
12
13     switch (whence)
14     {
15     case LSEEK_SET:
16         break;
17     case LSEEK_CUR:
18         pos += file->pos;
19         break;
20     case LSEEK_END:
21         if ((ret = vop_fstat(file->node, stat)) == 0)
22         {
23             pos += stat->st_size;
24         }
25         break;
26     default:
27         ret = -EINVAL;
28     }
29
30     if (ret == 0)
31     {
32         if ((ret = vop_tryseek(file->node, pos)) == 0)
33         {
34             file->pos = pos;
35         }
36         // cprintf("file_seek, pos=%d, whence=%d, ret=%d\n", pos,
whence, ret);
37     }
38     fd_array_release(file);
39     return ret;
40 }

```

该函数接收三个参数：文件描述符 `fd`，偏移量 `pos` 和定位的起始位置 `whence`，其中 `whence` 有三个选项：

- 文件开头（`LSEEK_SET`）
- 当前位置（`LSEEK_CUR`）
- 文件结尾（`LSEEK_END`）

该函数的作用是更改文件操作的起始点

2. 使用 `sysfile_read` 函数从文件中读取长度为 `len` 的数据，将其存储到 `buf` 指向的内存缓冲区中。

- 如果读取的数据长度不等于 `len`，则返回一个相应的错误代码
- 如果读取出现错误，返回 -1
- 如果读取成功，返回 0

`sysfile_read` 函数如下：

```

1  /* sysfile_read - read file */
2  int sysfile_read(int fd, void *base, size_t len)

```



```

3 {
4     struct mm_struct *mm = current->mm;
5     //检查读取长度是否为0
6     if (len == 0)
7     {
8         return 0;
9     }
10    //检查文件是否可读
11    if (!file_testfd(fd, 1, 0))
12    {
13        return -E_INVAL;
14    }
15    //分配buffer空间，调用kmalloc函数分配4096字节的buffer空间
16    void *buffer;
17    if ((buffer = kmalloc(IOBUF_SIZE)) == NULL)
18    {
19        return -E_NO_MEM;
20    }
21
22    int ret = 0;
23    size_t copied = 0, alen;
24    //循环读取文件
25    while (len != 0)
26    {
27        //每次读取buffer大小
28        //如果剩余部分小于4096字节，则只读取剩余部分
29        if ((alen = IOBUF_SIZE) > len)
30        {
31            alen = len;
32        }
33        //调用file_read函数将文件内容读取到buffer中，alen为实际大小
34        ret = file_read(fd, buffer, alen, &alen);
35        if (alen != 0)
36        {
37            lock_mm(mm);
38            {
39                //调用copy_to_user函数将读取到的内容拷贝到用户的内存空间中
40                if (copy_to_user(mm, base, buffer, alen))
41                {
42                    assert(len >= alen);
43                    base += alen, len -= alen, copied += alen;
44                }
45                else if (ret == 0)
46                {
47                    ret = -E_INVAL;
48                }
49            }
50            unlock_mm(mm);
51        }
52        if (ret != 0 || alen == 0)
53        {
54            goto out;

```

```

55     }
56 }
57 //函数调用层返回至用户程序
58 out:
59     kfree(buffer);
60     if (copied != 0)
61     {
62         return copied;
63     }
64     return ret;
65 }

```

这个函数的作用是从指定文件描述符对应的文件中读取数据到用户提供的缓存区 `base` 中

于是可以得到 `load_icode_read` 函数的作用是读取程序文件的原始数据内容：从指定文件描述符 `fd` 对应的文件中从 `offset` 偏移量开始读取 `len` 长度的字节到 `buf` 中。

## mm\_map()

函数原型如下：

```

1  int
2  mm_map(struct mm_struct *mm, uintptr_t addr, size_t len, uint32_t vm_flags,
3         struct vma_struct **vma_store) {
4      uintptr_t start = ROUNDDOWN(addr, PGSIZE), end = ROUNDUP(addr + len,
5      PGSIZE);
6      if (!USER_ACCESS(start, end)) {
7          return -E_INVAL;
8      }
9      assert(mm != NULL);
10
11     int ret = -E_INVAL;
12
13     struct vma_struct *vma;
14     if ((vma = find_vma(mm, start)) != NULL && end > vma->vm_start) {
15         goto out;
16     }
17     ret = -E_NO_MEM;
18
19     if ((vma = vma_create(start, end, vm_flags)) == NULL) {
20         goto out;
21     }
22     insert_vma_struct(mm, vma);
23     if (vma_store != NULL) {
24         *vma_store = vma;
25     }
26     ret = 0;
27
28 out:
29     return ret;
30 }

```

具体操作如下：

1. 确保地址范围与页对齐：根据给定的地址和长度，对地址进行向下取整（`ROUNDDOWN`）并对长度进行向上取整（`ROUNDUP`）
2. 检查所提供的地址范围是否属于用户空间
3. 查找 `mm` 中的 `vma` 是否与新的 `vma` 区域重叠
4. 调用 `vma_create` 创建一个新的 `vma`
5. 将新 `vma` 插入到进程的 `vma` 列表中
6. 将新的 `vma` 结构体存储到 `vma_store` 指向的位置

于是可以得到 `mm_map` 函数的作用是根据 `addr` 和 `len` 创建新的虚拟内存区域，并将其加入到 `mm` 的虚拟内存列表中，同时存储到 `vma_store` 指向的位置。

### `pgdir_alloc_page()`

函数原型如下：

```
1 //kern/mm/pmm.c
2
3 // pgdir_alloc_page - call alloc_page & page_insert functions to
4 //                      - allocate a page size memory & setup an addr map
5 //                      - pa<->la with linear address la and the PDT pgdir
6 struct Page *pgdir_alloc_page(pde_t *pgdir, uintptr_t la, uint32_t perm) {
7     struct Page *page = alloc_page();
8     if (page != NULL) {
9         if (page_insert(pgdir, page, la, perm) != 0) {
10             free_page(page);
11             return NULL;
12         }
13         if (swap_init_ok) {
14             if (check_mm_struct != NULL) {
15                 swap_map_swappable(check_mm_struct, la, page, 0);
16                 page->pra_vaddr = la;
17                 assert(page_ref(page) == 1);
18                 // cprintf("get No. %d page: pra_vaddr %x, pra_link.prev %x,
19                 // pra_link_next %x in pgdir_alloc_page\n", (page-pages),
20                 // page->pra_vaddr, page->pra_page_link.prev,
21                 // page->pra_page_link.next);
22             } else { // now current is existed, should fix it in the future
23                 // swap_map_swappable(current->mm, la, page, 0);
24                 // page->pra_vaddr=la;
25                 // assert(page_ref(page) == 1);
26                 // panic("pgdir_alloc_page: no pages. now current is existed,
27                 // should fix it in the future\n");
28             }
29         }
30     }
31
32     return page;
33 }
```

该函数作用是为给定的页目录表 *pgdir* 分配一个物理页面，页面的访问权限是 *perm*，并在给定的线性地址 *la* 处设置地址映射关系。

### (3.1) read raw data content in file and resolve elfhdr

代码如下：

```
1 struct Page *page;
2
3 struct elfhdr __elf, *elf = &__elf;
4 if ((ret = load_icode_read(fd, elf, sizeof(struct elfhdr), 0)) != 0)
5 {
6     goto bad_elf_cleanup_pgdir;
7 }
8
9 // 通过elf魔数判断文件是否有效
10 if (elf->e_magic != ELF_MAGIC)
11 {
12     ret = -E_INVAL_ELF;
13     goto bad_elf_cleanup_pgdir;
14 }
```

逻辑如下：

1. 创建一个指向 `elfhdr` 结构体的指针 `elf`
2. 调用 `load_icode_read` 函数将文件开头 `elfhdr` 结构大小的数据赋值给该指针
3. 通过 `elf->e_magic` 判断该文件是否有效
  - `elf->e_magic==ELF_MAGIC` 则说明该文件有效，没有损坏
  - 反之进入 `bad_elf_cleanup_pgdir` 部分

### (3.2) read raw data content in file and resolve proghdr based on info in elfhdr

代码如下：

```
1 struct proghdr __ph, *ph = &__ph;
2 uint32_t vm_flags, perm, phnum;
3 for (phnum = 0; phnum < elf->e_phnum; phnum++)
4 {
5     off_t phoff = elf->e_phoff + sizeof(struct proghdr) * phnum;
6     if ((ret = load_icode_read(fd, ph, sizeof(struct proghdr), phoff)) !=
7         0)
8     {
9         goto bad_cleanup_mmap;
10    }
11    // find every program section headers
12    if (ph->p_type != ELF_PT_LOAD)
13    {
14        continue;
15    }
16 }
```

```

14     }
15     if (ph->p_filesz > ph->p_memsz)
16     {
17         ret = -E_INVALID_FILE;
18         goto bad_cleanup_mmap;
19     }
20     if (ph->p_filesz == 0)
21     {
22         // continue ;
23     }

```

逻辑如下：

1. 创建一个指向 `proghdr` 结构体的指针 `ph`
2. `phnum` 用于遍历 `elf->e_phnum`（程序头表中的条目数）
3. `phoff` 表示每一个条目中文件偏移量：
  - 每一个条目所占的空间是 `struct proghdr` 大小
  - 本条目前面总共有 `phnum` 个条目
  - 再加上程序头表的文件偏移量 `elf->e_phoff` 即可得到
4. 调用 `load_icode_read` 函数将文件开头偏移量 `phoff` 开始的 `proghdr` 结构大小的数据赋值给该指针
5. 进行一些检查：
  - 程序段的类型是否是可加载的
  - 段在 ELF 文件中占据的空间大小是否不大于段加载到内存后占据的空间大小
  - 段在 ELF 文件中占据的空间大小是否不等于0

### (3.3) call `mm_map` to build vma related to TEXT/DATA

代码如下：

```

1     uint32_t vm_flags=0, perm=PTE_U | PTE_V;
2
3     if (ph->p_flags & ELF_PF_X)
4         vm_flags |= VM_EXEC;
5     if (ph->p_flags & ELF_PF_W)
6         vm_flags |= VM_WRITE;
7     if (ph->p_flags & ELF_PF_R)
8         vm_flags |= VM_READ;
9     // modify the perm bits here for RISC-V
10    if (vm_flags & VM_READ)
11        perm |= PTE_R;
12    if (vm_flags & VM_WRITE)
13        perm |= (PTE_W | PTE_R);
14    if (vm_flags & VM_EXEC)
15        perm |= PTE_X;
16    if ((ret = mm_map(mm, ph->p_va, ph->p_memsz, vm_flags, NULL)) != 0)
17    {
18        goto bad_cleanup_mmap;

```

逻辑如下：

1. `vm_flags`：表示虚拟内存区域的标志位（flags），初始化为0，表示还没有赋予权限

- 权限定义如下：

```
1 //kern/mm/vmm.h
2
3 #define VM_READ          0x00000001
4 #define VM_WRITE         0x00000002
5 #define VM_EXEC          0x00000004
6 #define VM_STACK         0x00000008
```

`perm`：表示页面权限，初始化为 `PTE_U` 和 `PTE_V`，表示用户权限和可执行权限

- 权限定义如下：

```
1 //kern/mm/mmu.h
2
3 // page table entry (PTE) fields
4 #define PTE_V          0x001 // valid
5 #define PTE_R          0x002 // Read
6 #define PTE_W          0x004 // Write
7 #define PTE_X          0x008 // Execute
8 #define PTE_U          0x010 // User
9 #define PTE_G          0x020 // Global
10 #define PTE_A          0x040 // Accessed
11 #define PTE_D          0x080 // Dirty
12 #define PTE_SOFT       0x300 // Reserved for Software
13
14 #define PAGE_TABLE_DIR (PTE_V)
15 #define READ_ONLY (PTE_R | PTE_V)
16 #define READ_WRITE (PTE_R | PTE_W | PTE_V)
17 #define EXEC_ONLY (PTE_X | PTE_V)
18 #define READ_EXEC (PTE_R | PTE_X | PTE_V)
19 #define READ_WRITE_EXEC (PTE_R | PTE_W | PTE_X | PTE_V)
20
21 #define PTE_USER (PTE_R | PTE_W | PTE_X | PTE_U | PTE_V)
```

2.
  - 如果段属性是可执行的（`ELF_PF_X`），虚拟内存区域的标志位也加上可执行（`VM_EXEC`）
  - 如果段属性是可写的（`ELF_PF_W`），虚拟内存区域的标志位也加上可写（`VM_WRITE`）
  - 如果段属性是可读的（`ELF_PF_R`），虚拟内存区域的标志位也加上可读（`VM_READ`）
3.
  - 如果该虚拟内存区域是可读的（`VM_READ`），则页面权限也加上可读（`PTE_R`）
  - 如果该虚拟内存区域是可写的（`VM_WRITE`），则页面权限也加上可写（`PTE_W`）
  - 如果该虚拟内存区域是可执行的（`VM_EXEC`），则页面权限也加上可执行（`PTE_X`）

- 调用 `mm_map` 函数根据 `ph->p_va` (段在虚拟地址空间中的起始地址) 和 `ph->p_memsz` (段在内存中的大小) 创建新的虚拟内存区域并将其链入 `mm` 中的虚拟内存区域列表, 该新建的虚拟内存区域的属性存在 `vm_flags` 中

### (3.4) call `pgdir_alloc_page` to allocate page for TEXT/DATA, read contents in file and copy them into the new allocated pages

为数据段和代码段分配页。

代码如下:

```
1      off_t offset = ph->p_offset;
2      size_t off, size;
3      uintptr_t start = ph->p_va;
4      uintptr_t end = ph->p_va + ph->p_filesz;
5      uintptr_t la = ROUNDDOWN(start, PGSIZE);
6
7      ret = -E_NO_MEM;
8
9      while (start < end)
10     {
11         if ((page = pgdir_alloc_page(mm->pgdir, la, perm)) == NULL)
12         {
13             ret = -E_NO_MEM;
14             goto bad_cleanup_mmap;
15         }
16         off = start - la, size = PGSIZE - off, la += PGSIZE;
17         if (end < la)
18         {
19             size -= la - end;
20         }
21         if ((ret = load_icode_read(fd, page2kva(page) + off, size, offset))
22             != 0)
23         {
24             goto bad_cleanup_mmap;
25         }
26         start += size, offset += size;
27     }
```

逻辑如下:

- 首先初始化 `offset` 为程序段在文件中的偏移量, `start` 为加载的程序段的起始虚拟地址, `end` 表示加载的程序段的结束虚拟地址
- 接下来进行程序段的逐页加载:
  - 为 `pgdir` 分配一个物理页, 页面权限是 `perm`
  - 计算当前页内的偏移量 `off` 和剩余空间大小 `size`
  - 如果该页的大小超出程序段的实际大小, 则减去多余的部分
  - 调用 `load_icode_read` 函数从 `fd` 表示的文件中偏移量为 `offset` 的地址处读取 `size` 大小的数据到 `page2kva(page) + off` 中



- `page2kva(page)` 用于将给定的物理页面 `page` 转换为内核虚拟地址
  - `off` 是相对于这个物理页面起始地址的偏移量
  - 虽然起始地址不同但是两个偏移量是相同的，所以两者相加可以得到这个物理页面上对应偏移量处的内核虚拟地址
- 更新当前加载的位置和文件偏移量

### (3.5) call `pgdir_alloc_page` to allocate pages for BSS, memset zero in these pages

为 BSS 段分配页。

代码如下：

```
1      end = ph->p_va + ph->p_memsz;
2      if (start < la)
3      {
4          /* ph->p_memsz == ph->p_filesz */
5          if (start == end)
6          {
7              continue;
8          }
9          off = start + PGSIZE - la, size = PGSIZE - off;
10         if (end < la)
11         {
12             size -= la - end;
13         }
14         memset(page2kva(page) + off, 0, size);
15         start += size;
16         assert((end < la && start == end) || (end >= la && start == la));
17     }
18     while (start < end)
19     {
20         if ((page = pgdir_alloc_page(mm->pgdir, la, perm)) == NULL)
21         {
22             ret = -E_NO_MEM;
23             goto bad_cleanup_mmap;
24         }
25         off = start - la, size = PGSIZE - off, la += PGSIZE;
26         if (end < la)
27         {
28             size -= la - end;
29         }
30         memset(page2kva(page) + off, 0, size);
31         start += size;
32     }
33 }
34 sysfile_close(fd);
```

逻辑如下：

1. `end` 表示程序段在内存中结束的地址
2. `start` 表示当前处理的物理页面的起始地址，if段进行了段对齐的操作：如果 `start` 小于 `1a`，即条件成立，则需要对未被覆盖到的内存区域填充0
3. 使用循环逐个分配物理页面，将程序段的内容加载到内存中
  - 调用 `pgdir_alloc_page` 函数分配一个物理页面，页面权限是 `perm`
  - 计算偏移量 `off` 和大小 `size`
  - 判断是否有未被覆盖的内存区域，将该部分填充为0
4. 调用 `sysfile_close` 函数关闭文件

## (4) call `mm_map` to setup user stack, and put parameters into user stack

设置用户栈。

代码如下：

```
1   vm_flags = VM_READ | VM_WRITE | VM_STACK;
2   if ((ret = mm_map(mm, USTACKTOP - USTACKSIZE, USTACKSIZE, vm_flags, NULL)) !=
0)
3   {
4       goto bad_cleanup_mmap;
5   }
6   assert(pgdir_alloc_page(mm->pgdir, USTACKTOP - PGSIZE, PTE_USER) != NULL);
7   assert(pgdir_alloc_page(mm->pgdir, USTACKTOP - 2 * PGSIZE, PTE_USER) !=
NULL);
8   assert(pgdir_alloc_page(mm->pgdir, USTACKTOP - 3 * PGSIZE, PTE_USER) !=
NULL);
9   assert(pgdir_alloc_page(mm->pgdir, USTACKTOP - 4 * PGSIZE, PTE_USER) !=
NULL);
```

逻辑如下：

1. 初始化虚拟内存空间访问权限为可读可写并且该内存区域用作栈
2. 调用 `mm_map` 函数设置用户栈
3. 调用 `pgdir_alloc_page` 函数在用户栈顶的不同位置分配了多个页面，并使用 `assert` 断言确保这些页面的分配是成功的，每个页面的权限标志均为 `PTE_USER`，实现了将参数放入用户栈的功能

## (5) setup current process's mm, cr3, reset pgdir (using `lcr3` MARCO)

代码如下：

```
1   mm_count_inc(mm);
2   current->mm = mm;
3   current->cr3 = PADDR(mm->pgdir);
4   lcr3(PADDR(mm->pgdir));
```

逻辑如下：

1. 增加了 `mm` 的引用计数
2. 将当前进程的地址空间指针设置为 `mm`
3. 将当前目录基址寄存器 `cr3` 的值设置为当前进程页目录表的物理地址
4. 通过 `1cr3` 汇编指令实现页目录表的加载

## (6) setup uargc and uargv in user stacks

代码如下：

```
1  uint32_t argv_size = 0, i;
2  for (i = 0; i < argc; i++)
3  {
4      argv_size += strlen(kargv[i], EXEC_MAX_ARG_LEN + 1) + 1;
5  }
6
7  uintptr_t stacktop = USTACKTOP - (argv_size / sizeof(long) + 1) *
sizeof(long);
8  char **uargv = (char **)(stacktop - argc * sizeof(char *));
9
10 argv_size = 0;
11 for (i = 0; i < argc; i++)
12 {
13     uargv[i] = strcpy((char *)(stacktop + argv_size), kargv[i]);
14     argv_size += strlen(kargv[i], EXEC_MAX_ARG_LEN + 1) + 1;
15 }
16
17 stacktop = (uintptr_t)uargv - sizeof(int);
18 *(int *)stacktop = argc;
```

逻辑如下：

1. `argv_size` 初始化为0，用来存储所有参数的总大小；`i`用于遍历
2. 通过一个循环，计算所有参数的总大小，存在 `argv_size` 里面
  - `kargv` 表示参数列表，`kargv[i]` 表示每一个参数
  - `strlen` 函数原型是 `size_t strlen(const char *s, size_t maxlen)`；有一个最大参数限制 `maxlen`，返回扫描到的 `s` 的实际长度
  - 最后算出的长度值还要加1的原因是还要加上休止字符串 `'\0'`
3. 根据 `argv_size`，计算出栈顶位置 `stacktop`
  - `argv_size / sizeof(long)` 计算出参数列表所占用的长整型数目
  - `+1` 是因为还有一个 `NULL` 指针
  - `* sizeof(long)` 将长整型的数目转换回字节，获得参数列表的总字节数
4. 为参数数组 `uargv` 分配空间并设置其位置
  - `argc * sizeof(char *)` 计算出参数指针数组所占用的字节数

- 将参数指针数组的字节数从用户栈的顶部减去，即为指向参数指针数组的指针 `uargv` 的位置
- 5. 清空 `argv_size`，该变量将表示在程序执行时获取的参数个数
- 6. 使用循环，对每一个参数字符串做变动：
  - 调用 `strcpy` 函数将第 `i` 个参数字符串复制到对应的用户栈空间（`stacktop + argv_size`）上
  - 更新 `argv_size`，表示接下来的参数应该放于哪个位置
- 7. 将 `argc` 存储在用户栈的最顶端

## (7) setup trapframe for user environment

代码如下：

```

1  struct trapframe *tf = current->tf;
2  // Keep sstatus
3  uintptr_t sstatus = tf->status;
4  memset(tf, 0, sizeof(struct trapframe));
5  // Set the user stack top
6  tf->gpr.sp = stacktop;
7  // Set the entry point of the user program
8  tf->epc = elf->e_entry;
9  // Set the status register for the user program
10 tf->status = sstatus & ~(SSTATUS_SPP | SSTATUS_SPIE);
11 ret = 0;

```

逻辑如下：

1. 获取当前进程控制块中中断帧的指针
2. 保存中断帧中的 `status` 中的值于 `sstatus` 中
3. 清空中断帧的内容
4. 设置用户栈顶指针
5. 设置用户程序的入口点
6. 恢复保存的 `status` 的值同时清除 `SSTATUS_SPP` 和 `SSTATUS_SPIE` 标志位

## 完整代码

```

1  static int
2  load_icode(int fd, int argc, char **kargv)
3  {
4      /* LAB8:EXERCISE2 YOUR CODE HINT:2113285
5       * how to load the file with handler fd in to process's memory? how to
6       * setup argc/argv?
7       * MACROS or Functions:
8       * mm_create      - create a mm
9       * setup_pgdir    - setup pgdir in mm
10      * load_icode_read - read raw data content of program file
11      * mm_map         - build new vma
12      * pgdir_alloc_page - allocate new memory for TEXT/DATA/BSS/stack parts

```

```

12      * lcr3          - update Page Directory Addr Register -- CR3
13      */
14      // You can Follow the code form LAB5 which you have completed to complete
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
19      * (3.2) read raw data content in file and resolve proghdr based on
info in elfhdr
20      * (3.3) call mm_map to build vma related to TEXT/DATA
21      * (3.4) call pgdir_alloc_page to allocate page for TEXT/DATA, read
contents in file
22      *          and copy them into the new allocated pages
23      * (3.5) call pgdir_alloc_page to allocate pages for BSS, memset zero
in these pages
24      * (4) call mm_map to setup user stack, and put parameters into user stack
25      * (5) setup current process's mm, cr3, reset pgidr (using lcr3 MARCO)
26      * (6) setup uargc and uargv in user stacks
27      * (7) setup trapframe for user environment
28      * (8) if up steps failed, you should cleanup the env.
29      */
30
31      // 当前内存空间必须为空
32      if (current->mm != NULL)
33      {
34          panic("load_icode: current->mm must be empty.\n");
35      }
36
37      int ret = -E_NO_MEM;
38      struct mm_struct *mm;
39
40      //(1) create a new mm for current process
41      if ((mm = mm_create()) == NULL)
42      {
43          goto bad_mm;
44      }
45
46      //(2) create a new PDT, and mm->pgdir= kernel virtual addr of PDT
47      if (setup_pgdir(mm) != 0)
48      {
49          goto bad_pgdir_cleanup_mm;
50      }
51
52      //(3) copy TEXT/DATA/BSS parts in binary to memory space of process
53      struct Page *page;
54
55      //(3.1) read raw data content in file and resolve elfhdr
56      struct elfhdr __elf, *elf = &__elf;
57      if ((ret = load_icode_read(fd, elf, sizeof(struct elfhdr), 0)) != 0)
58      {
59          goto bad_elf_cleanup_pgdir;
60      }

```

```

61
62 // 通过elf判断文件是否有效
63 if (elf->e_magic != ELF_MAGIC)
64 {
65     ret = -E_INVALID_ELF;
66     goto bad_elf_cleanup_pgdir;
67 }
68
69 //(3.2) read raw data content in file and resolve proghdr based on info in
elfhdr
70 struct proghdr __ph, *ph = &__ph;
71 uint32_t phnum;
72 for (phnum = 0; phnum < elf->e_phnum; phnum++)
73 {
74     off_t phoff = elf->e_phoff + sizeof(struct proghdr) * phnum;
75     if ((ret = load_icode_read(fd, ph, sizeof(struct proghdr), phoff)) !=
0)
76     {
77         goto bad_cleanup_mmap;
78     }
79     // find every program section headers
80     if (ph->p_type != ELF_PT_LOAD)
81     {
82         continue;
83     }
84     if (ph->p_filesz > ph->p_memsz)
85     {
86         ret = -E_INVALID_ELF;
87         goto bad_cleanup_mmap;
88     }
89     if (ph->p_filesz == 0)
90     {
91         // continue ;
92     }
93
94     //(3.3) call mm_map to build vma related to TEXT/DATA
95     uint32_t vm_flags=0, perm=PTE_U | PTE_V;
96     if (ph->p_flags & ELF_PF_X)
97         vm_flags |= VM_EXEC;
98     if (ph->p_flags & ELF_PF_W)
99         vm_flags |= VM_WRITE;
100     if (ph->p_flags & ELF_PF_R)
101         vm_flags |= VM_READ;
102     // modify the perm bits here for RISC-V
103     if (vm_flags & VM_READ)
104         perm |= PTE_R;
105     if (vm_flags & VM_WRITE)
106         perm |= (PTE_W | PTE_R);
107     if (vm_flags & VM_EXEC)
108         perm |= PTE_X;
109     if ((ret = mm_map(mm, ph->p_va, ph->p_memsz, vm_flags, NULL)) != 0)
110     {

```

```

111         goto bad_cleanup_mmap;
112     }
113
114     //(3.4) call pgdir_alloc_page to allocate page for TEXT/DATA, read
contents in file
115     // and copy them into the new allocated pages
116     off_t offset = ph->p_offset;
117     size_t off, size;
118     uintptr_t start = ph->p_va;
119     uintptr_t end= ph->p_va + ph->p_filesz;
120     uintptr_t la = ROUNDDOWN(start, PGSIZE);
121
122     ret = -E_NO_MEM;
123
124     while (start < end)
125     {
126         if ((page = pgdir_alloc_page(mm->pgdir, la, perm)) == NULL)
127         {
128             ret = -E_NO_MEM;
129             goto bad_cleanup_mmap;
130         }
131         off = start - la, size = PGSIZE - off, la += PGSIZE;
132         if (end < la)
133         {
134             size -= la - end;
135         }
136         if ((ret = load_icode_read(fd, page2kva(page) + off, size,
offset)) != 0)
137         {
138             goto bad_cleanup_mmap;
139         }
140         start += size, offset += size;
141     }
142
143     //(3.5) call pgdir_alloc_page to allocate pages for BSS, memset zero
in these pages
144     end = ph->p_va + ph->p_memsz;
145     if (start < la)
146     {
147         /* ph->p_memsz == ph->p_filesz */
148         if (start == end)
149         {
150             continue;
151         }
152         off = start + PGSIZE - la, size = PGSIZE - off;
153         if (end < la)
154         {
155             size -= la - end;
156         }
157         memset(page2kva(page) + off, 0, size);
158         start += size;
159         assert((end < la && start == end) || (end >= la && start == la));

```



```

160     }
161     while (start < end)
162     {
163         if ((page = pgdir_alloc_page(mm->pgdir, la, perm)) == NULL)
164         {
165             ret = -E_NO_MEM;
166             goto bad_cleanup_mmap;
167         }
168         off = start - la, size = PGSIZE - off, la += PGSIZE;
169         if (end < la)
170         {
171             size -= la - end;
172         }
173         memset(page2kva(page) + off, 0, size);
174         start += size;
175     }
176 }
177 sysfile_close(fd);
178
179 // (4) call mm_map to setup user stack, and put parameters into user stack
180 vm_flags = VM_READ | VM_WRITE | VM_STACK;
181 if ((ret = mm_map(mm, USTACKTOP - USTACKSIZE, USTACKSIZE, vm_flags, NULL))
182 != 0)
183 {
184     goto bad_cleanup_mmap;
185 }
186 assert(pgdir_alloc_page(mm->pgdir, USTACKTOP - PGSIZE, PTE_USER) != NULL);
187 assert(pgdir_alloc_page(mm->pgdir, USTACKTOP - 2 * PGSIZE, PTE_USER) !=
188 NULL);
189 assert(pgdir_alloc_page(mm->pgdir, USTACKTOP - 3 * PGSIZE, PTE_USER) !=
190 NULL);
191 assert(pgdir_alloc_page(mm->pgdir, USTACKTOP - 4 * PGSIZE, PTE_USER) !=
192 NULL);
193
194 // (5) setup current process's mm, cr3, reset pgdir (using lcr3 MARCO)
195 mm_count_inc(mm);
196 current->mm = mm;
197 current->cr3 = PADDR(mm->pgdir);
198 lcr3(PADDR(mm->pgdir));
199
200 // (6) setup uargc and uargv in user stacks
201 uint32_t argv_size = 0, i;
202 for (i = 0; i < argc; i++)
203 {
204     argv_size += strlen(kargv[i], EXEC_MAX_ARG_LEN + 1) + 1;
205 }
206
207 uintptr_t stacktop = USTACKTOP - (argv_size / sizeof(long) + 1) *
208 sizeof(long);
209 char **uargv = (char **)(stacktop - argc * sizeof(char *));
210
211 argv_size = 0;

```

```

207     for (i = 0; i < argc; i++)
208     {
209         uargv[i] = strcpy((char *) (stacktop + argv_size), kargv[i]);
210         argv_size += strlen(kargv[i], EXEC_MAX_ARG_LEN + 1) + 1;
211     }
212
213     stacktop = (uintptr_t)uargv - sizeof(int);
214     *(int *)stacktop = argc;
215
216     //(7) setup trapframe for user environment
217     struct trapframe *tf = current->tf;
218     // Keep sstatus
219     uintptr_t sstatus = tf->sstatus;
220     memset(tf, 0, sizeof(struct trapframe));
221     // Set the user stack top
222     tf->gpr.sp = stacktop;
223     // Set the entry point of the user program
224     tf->epc = elf->e_entry;
225     // Set the status register for the user program
226     tf->sstatus = sstatus & ~(SSTATUS_SPP | SSTATUS_SPIE);
227     ret = 0;
228 out:
229     return ret;
230 bad_cleanup_mmap:
231     exit_mmap(mm);
232 bad_elf_cleanup_pgdir:
233     put_pgdir(mm);
234 bad_pgdir_cleanup_mm:
235     mm_destroy(mm);
236 bad_mm:
237     goto out;
238 }

```

## 结果展示

- 输入命令 `hello`

```

Hello world!!.
I am process 3.
hello pass.

```

- 输入命令 `forktest`

```
I am child 31
I am child 30
I am child 29
I am child 28
I am child 27
I am child 26
I am child 25
I am child 24
I am child 23
I am child 22
I am child 21
I am child 20
I am child 19
I am child 18
I am child 17
I am child 16
I am child 15
I am child 14
I am child 13
I am child 12
I am child 11
I am child 10
I am child 9
I am child 8
I am child 7
I am child 6
I am child 5
I am child 4
I am child 3
I am child 2
I am child 1
I am child 0
forktest pass.
```

- 输入命令 `matrix`

```
pid 51 is running (3500 times)!.  
pid 51 done!.  
pid 53 is running (5900 times)!.  
pid 53 done!.  
pid 55 is running (15400 times)!.  
pid 58 is running (17900 times)!.  
pid 58 done!.  
pid 47 done!.  
pid 55 done!.  
pid 43 done!.  
pid 44 done!.  
pid 50 done!.  
pid 52 done!.  
pid 57 done!.  
pid 45 done!.  
pid 54 done!.  
pid 48 done!.  
matrix pass.
```

- 输入命令 `yield`

```
Hello, I am process 3.  
Back in process 3, iteration 0.  
Back in process 3, iteration 1.  
Back in process 3, iteration 2.  
Back in process 3, iteration 3.  
Back in process 3, iteration 4.  
All done in process 3.  
yield pass.
```

- 输入命令 `testbss`

```
Making sure bss works right...  
Yes, good. Now doing a wild write off the end...  
testbss may pass.
```

**CH1: 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**

---

## 进程间通信

每个进程各自有不同的用户地址空间，任何一个进程的全局变量在另一个进程中都看不到，所以进程之间要交换数据必须通过内核，在内核中开辟一块缓冲区，进程A把数据从用户空间拷到内核缓冲区，进程B再从内核缓冲区把数据读走，内核提供的这种机制称为进程间通信。

## 管道

管道是一种最基本的进程间通信机制，把一个进程连接到另一个进程的一个数据流称为一个管道，通常是用作**把一个进程的输出通过管道连接到另一个进程的输入**

管道本质上是内核的一块缓存，内核维护了一块缓冲区与管道文件相关联，对管道文件的操作被内核转换成对这块缓冲区内内存的操作

## 设计思路

1. 考虑在磁盘上保留一部分空间来作为pipe机制的缓冲区；
2. 同时当要求建立管道时，通过两个进程的进程控制块来标识，表示输出数据时需要输出到管道的输出缓冲区中而读取数据时要从输入缓冲区中读取数据。

## 定义管道数据结构

在内核中定义一个管道数据结构来存储进程间传输的数据。该数据结构包括两个缓冲区，一个输入缓冲区和一个输出缓冲区，每个缓冲区由一个首指针和一个尾指针组成，以便于进行读写操作。

## 管道创建与销毁

管道创建：当进程A需要向进程B发送数据时，他调用pipe()的系统调用，该系统调用将创建一个新的管道，并返回一个标识符。进程A使用该标识符写入数据到特定管道的输出缓冲区中，进程B使用特定标识符从相同的管道中的输入缓冲区中读取数据。

管道销毁：当不再需要管道时，进程调用close()系统调用来关闭读或写端口，当所有的可以使用该管道的进程都已经不需要该管道时，管道销毁

## 管道读写操作

管道读操作：使用read()系统调用从输入缓冲区中读取数据

管道写操作：使用write()系统调用向输出缓冲区中写数据

**注意：如果输入缓冲区为空，则读操作会不被允许执行，直到有某个进程向其中写数据；如果缓冲区已满，则写操作不被允许，直到有某个数据被进程读出并且有足够的空间供该进程写入为止。**

## 管道进程同步

因为管道中的数据是共享的，因此可能需要使用信号量或者互斥锁等机制来保证多个进程之间的同步：

- 写入数据前应该加锁，读取数据时也应该加锁
- 使用信号量，保证写操作和读操作的先后性其他限制：
  - 管道中数据为空时进行读操作时，读操作的进程进入等待队列
  - 管道中数据满时进行写操作时，写操作的进程需要等待

## 虚拟文件系统的实现

为了兼容UNIX的标准I/O接口，可以在虚拟我呢见系统中实现管道机制，以方便管道使用标准I/O接口进行读写操作

## CH2：完成基于“UNIX的软连接和硬连接机制”的设计方案

---

### UNIX的硬链接和软链接机制

#### 硬链接

通过具有相同 `inode` 号的链接相互连接的文件

特点：

- 文件有相同的 `inode` 及 `data block`；
- 只能对已存在的文件进行创建；
- 不能交叉文件系统进行硬链接的创建；
- 不能对目录进行创建，只可对文件创建；
- 删除一个硬链接文件并不影响其他有相同 `inode` 号的文件

优点：删除、重命名或将目标文件移动到新位置不会使硬链接变得毫无价值。它在保持不同 `inode` 的同时变得更加高效和实用

缺点：难以处理

#### 软链接

又叫符号链接，通过具有不同索引节点号的链接互连的文件

特点：

- 软链接有自己的文件属性及权限等；
- 可对不存在的文件或目录创建软链接；
- 软链接可交叉文件系统；
- 软链接可对文件或目录创建；
- 创建软链接时，链接计数 `i_nlink` 不会增加；
- 删除软链接并不影响被指向的文件，但若被指向的原文件被删除，则相关软连接被称为死链接

优点：每个文件，文件符号对象都可以与一个软链接相关联

缺点：删除目标文件或将目标文件移动到新位置会使软链接变得毫无价值

#### 总结

1. 软链接或符号链接通过路径指向文件或目录，而硬链接则指向磁盘上的数据。
2. 删除目标文件不会影响硬链接，但会使软链接失效。
3. 软链接可以链接到跨不同文件系统的文件或目录，而硬链接则不能。

# 设计方案

由于 `sfs_disk_inode` 结构体存在成员变量 `nlinks` 表示当前文件被链接的计数，所以支持硬链接和软链接机制。

## 硬链接

如果要创建文件A的硬链接B，首先将B当成正常文件创建

- 将B的TYPE域设置为链接，再使用一个位标记该链接是硬链接
- 将B的 `inode` 指向目标文件，即文件A的 `inode` → 硬链接的 `inode` 和目标文件的 `inode` 是一个
- 将文件A被链接的计数加1

访问到B时，判断B是一个链接，则实际访问的文件是B指向的文件A

删除链接B时，需要做两个操作：

- 删除B的 `inode`
- 将B指向的文件A的被链接计数 `nlinks` 减1，如果减到了0，则将A删除

## 软链接

如果要创建文件A的软链接B，首先将B当成正常文件创建

- 将B的TYPE域设置为链接，再使用一个位标记该链接是软链接
- `inode` 与文件A的 `inode` 不同
- 将文件A的地址存放到文件B的内容中去
- 被链接计数不变

访问到B的时候，判断B是一个链接，则实际访问的文件是B指向的文件A

当删除链接B时，直接将其在磁盘上的 `inode` 删掉即可

## 具体实现

### 用到的变量

- 1 `inode`: 表示文件的元数据，包括文件类型、大小、权限等信息。
- 2 `file`: 表示打开的文件，包括指向`inode`的指针、读写偏移量等信息。
- 3 `nlinks`: 硬链接中表示链接计数

### 接口信息

```
1 int link(const char *oldpath, const char *newpath);
2 //创建一个硬链接，将oldpath指向的文件与newpath指向的文件名关联起来
3 //在目标文件夹的控制块中增加一个描述符，并且两个描述符的inode指针相同，同时nlinks数据结构应该相应增加。
4
5 int symlink(const char *target, const char *linkpath);
6 //创建一个软链接，将target指向的地址与linkpath指向的文件名关联起来
7 //在文件目录中创建一个特殊的文件，其中包含指向target的地址信息。
```



```
8
9  int readlink(const char *pathname, char *buf, size_t bufsiz);
10 //读取软链接文件的内容，将软链接指向的地址信息放入buf中
11
12 int unlink(const char *pathname);
13 //删除一个链接
14 //如果是硬链接，减少相关nlinks数据结构的计数；
15 //如果是软链接，直接删除软链接文件
```