

# Lab9 file system

---

## Lab9 file system

前置知识

实验内容

Large files (moderate)

任务

理解文件系统结构

修改 `bmap` 及 `itrunc` 函数:

修改 `struct inode`

Symbolic links (moderate)

任务

步骤

实验得分

## 前置知识

## 实验内容

### Large files (moderate)

#### 任务

- 目前, xv6 文件系统 中的每个 `inode` 包含 12 个“直接”块号和一个“单一间接”块号, 后者引用一个块, 该块最多可以包含 256 个块号, 总共可以引用  $12+256=268$  个块。因此, xv6 文件的最大大小限制为 268 个块。

你需要修改 xv6 文件系统代码, 支持每个 `inode` 中的“双重间接”块, 这样文件的最大大小可以增加到 65803 个块, 即  $256*256+256+11$  个块 (因为我们将牺牲一个直接块号以用于双重间接块)。

## 理解文件系统结构

- `struct dinode`（定义在 `fs.h` 中）描述了磁盘上的 `inode` 结构。你需要关注 `NDIRECT`、`NINDIRECT`、`MAXFILE` 以及 `struct dinode` 中的 `addrs[]` 元素。
- `bmap()` 函数位于 `fs.c` 中，用于查找文件在磁盘上的数据块。`bmap()` 会根据需要分配新块以容纳文件内容，并在需要时分配间接块来存储块地址。

### 修改 `bmap` 及 `itrunc` 函数：

- 你需要在 `bmap()` 中实现双重间接块（doubly-indirect block），同时保持对直接块和单一间接块的支持。
- 将 `ip->addrs[]` 的前11个元素保留为直接块，第12个元素用作单一间接块，第13个元素用作新的双重间接块。
- 修改后的 `bmap()` 应该能够映射文件的逻辑块号（`bn`）到磁盘块号，并在需要时分配新的块。

```
1 // Return the disk block address of the nth block in inode ip.
2 // If there is no such block, bmap allocates one.
3 static uint
4 bmap(struct inode *ip, uint bn)
5 {
6     uint addr, *a;
7     struct buf *bp;
8
9     if(bn < NDIRECT){
10         if((addr = ip->addrs[bn]) == 0)
11             ip->addrs[bn] = addr = balloc(ip->dev);
12         return addr;
13     }
14     bn -= NDIRECT;
15
16     if(bn < NINDIRECT){
17         // Load indirect block, allocating if necessary.
18         if((addr = ip->addrs[NDIRECT]) == 0)
19             ip->addrs[NDIRECT] = addr = balloc(ip->dev);
20         bp = bread(ip->dev, addr);
21         a = (uint*)bp->data;
22         if((addr = a[bn]) == 0){
23             a[bn] = addr = balloc(ip->dev);
```

```

24     log_write(bp);
25 }
26 brelse(bp);
27 return addr;
28 }
29
30 bn -= NINDIRECT;
31 // 去除已经由直接块和单间接块映射的块数，以得到在双间接块中的相对块号
32
33 if (bn < NDBL_INDIRECT) {
34     // 如果文件的双间接块不存在，则分配一个
35     if ((addr = ip->addrs[NDIRECT + 1]) == 0) {
36         addr = balloc(ip->dev);
37         if (addr == 0)
38             return 0;
39         ip->addrs[NDIRECT + 1] = addr;
40     }
41
42     // 读取双间接块
43     bp = bread(ip->dev, addr);
44     a = (uint*)bp->data;
45
46     // 计算在单间接块数组中的索引，即第几个单间接块
47     uint index1 = bn / NINDIRECT;
48
49     // 如果这个单间接块不存在，则分配一个
50     if ((addr = a[index1]) == 0) {
51         addr = balloc(ip->dev);
52         if (addr == 0)
53             return 0;
54         a[bn / NINDIRECT] = addr;
55         log_write(bp); // Record changes in the log
56     }
57     brelse(bp);
58
59     // 读取相应的单间接块
60     bp = bread(ip->dev, addr);
61     a = (uint *)bp->data;
62
63     // 计算在单间接块中的索引，即单间接块中的第几个数据块
64     uint index2 = bn % NINDIRECT;
65

```

```

66 // 如果这个数据块不存在，则分配一个
67 if ((addr = a[index2]) == 0) {
68     addr = balloc(ip->dev);
69     if (addr == 0)
70         return 0;
71     a[bn % NINDIRECT] = addr;
72     log_write(bp); // Record changes in the log
73 }
74 brelse(bp);
75 return addr; // Returns the actual data block
76 }
77
78 panic("bmap: out of range");
79 }

```

```

1 // Truncate inode (discard contents).
2 // Caller must hold ip->lock.
3 void
4 itrunc(struct inode *ip)
5 {
6     int i, j;
7     struct buf *bp;
8     uint *a;
9
10    for(i = 0; i < NDIRECT; i++){
11        if(ip->addrs[i]){
12            bfree(ip->dev, ip->addrs[i]);
13            ip->addrs[i] = 0;
14        }
15    }
16
17    if(ip->addrs[NDIRECT]){
18        bp = bread(ip->dev, ip->addrs[NDIRECT]);
19        a = (uint*)bp->data;
20        for(j = 0; j < NINDIRECT; j++){
21            if(a[j])
22                bfree(ip->dev, a[j]);
23        }
24        brelse(bp);
25        bfree(ip->dev, ip->addrs[NDIRECT]);
26        ip->addrs[NDIRECT] = 0;

```

```

27     }
28
29     if (ip->addrs[NDIRECT + 1]) {
30         // 读取双间接块
31         bp = bread(ip->dev, ip->addrs[NDIRECT + 1]);
32         a = (uint*)bp->data;
33
34         for (i = 0; i < NINDIRECT; ++i) {
35             if (a[i] == 0) continue;
36
37             // 读取单间接块
38             struct buf* bp2 = bread(ip->dev, a[i]);
39             uint* b = (uint*)bp2->data;
40             for (j = 0; j < NINDIRECT; ++j) {
41                 if (b[j])
42                     bfree(ip->dev, b[j]); // 释放数据块
43             }
44             brelse(bp2);
45
46             bfree(ip->dev, a[i]); // 释放单间接块
47             a[i] = 0;
48         }
49         brelse(bp);
50
51         bfree(ip->dev, ip->addrs[NDIRECT + 1]); // 释放双间接块
52         ip->addrs[NDIRECT + 1] = 0;
53     }
54
55     ip->size = 0;
56     iupdate(ip);
57 }

```

## 修改 struct inode

- kernel/file.h:

```

1 //添加到结构体中
2 uint addrs[NDIRECT+2];

```

## Symbolic links (moderate)

### 任务

- 这个实验的目标是为xv6操作系统增加对符号链接（symbolic links）的支持。符号链接是一种特殊的文件类型，它通过路径名引用另一个文件。当打开符号链接时，内核会跟随链接指向的文件。符号链接类似于硬链接，但硬链接只能指向同一磁盘上的文件，而符号链接可以跨磁盘设备引用文件。虽然xv6不支持多个设备，但实现这个系统调用是理解路径名查找工作原理的好练习。
- 实现 `symlink` 系统调用：

**任务描述：** 你将实现 `symlink(char *target, char *path)` 系统调用，它在 `path` 位置创建一个新的符号链接，指向 `target` 文件。你需要确保在实现过程中符号链接能够正确地处理，尤其是处理路径名查找的递归和循环检测。

### 步骤

创建新系统调用：

- 为 `symlink` 创建一个新的系统调用编号，并在 `user/usys.pl` 和 `user/user.h` 中添加对应的条目。
- 在 `kernel/sysfile.c` 中实现一个空的 `sys_symlink` 函数。

```
1 // system calls
2 int fork(void);
3 int exit(int) __attribute__((noreturn));
4 int wait(int*);
5 int pipe(int*);
6 int write(int, const void*, int);
7 int read(int, void*, int);
8 int close(int);
9 int kill(int);
10 int exec(char*, char**);
11 int open(const char*, int);
12 int mknod(const char*, short, short);
13 int unlink(const char*);
14 int fstat(int fd, struct stat*);
15 int link(const char*, const char*);
16 int mkdir(const char*);
17 int chdir(const char*);
```

```
18 int dup(int);
19 int getpid(void);
20 char* sbrk(int);
21 int sleep(int);
22 int uptime(void);
23 //here
24 int symlink(char*, char*);
```

```
1 entry("fork");
2 entry("exit");
3 entry("wait");
4 entry("pipe");
5 entry("read");
6 entry("write");
7 entry("close");
8 entry("kill");
9 entry("exec");
10 entry("open");
11 entry("mknod");
12 entry("unlink");
13 entry("fstat");
14 entry("link");
15 entry("mkdir");
16 entry("chdir");
17 entry("dup");
18 entry("getpid");
19 entry("sbrk");
20 entry("sleep");
21 entry("uptime");
22 //here
23 entry("symlink");
```

定义符号链接文件类型：

- 在 `kernel/stat.h` 中添加一个新的文件类型 `T_SYMLINK`，用于表示符号链接。

```

1 #define T_DIR      1    // Directory
2 #define T_FILE     2    // File
3 #define T_DEVICE   3    // Device
4 #define T_SYMLINK  4    // Soft symbolic link - lab 9.2
5
6 struct stat {
7     int dev;        // File system's disk device
8     uint ino;       // Inode number
9     short type;     // Type of file
10    short nlink;    // Number of links to file
11    uint64 size;    // Size of file in bytes
12 };

```

添加新标志：

- 在 `kernel/fcntl.h` 中添加一个新的标志 `O_NOFOLLOW`，用于 `open` 系统调用。当指定该标志时，`open` 应该打开符号链接本身，而不是跟随符号链接指向的文件。确保该标志不会与现有标志冲突。

```

1 #define O_RDONLY  0x000
2 #define O_WRONLY  0x001
3 #define O_RDWR   0x002
4 #define O_CREATE  0x200
5 #define O_TRUNC   0x400
6 #define O_NOFOLLOW 0x004    // lab 9.2

```

实现 `symlink` 系统调用：

- 在 `path` 位置创建一个新的符号链接，链接到 `target`。符号链接的目标路径可以存储在 `inode` 的数据块中。
- `symlink` 应返回一个整数，表示成功（0）或失败（-1），类似于 `link` 和 `unlink` 系统调用。

```

1 // Create the path new as a link to the same inode as old.
2 uint64
3 sys_link(void)
4 {
5     char name[DIRSIZ], new[MAXPATH], old[MAXPATH];
6     struct inode *dp, *ip;
7

```



```
8  if(argstr(0, old, MAXPATH) < 0 || argstr(1, new, MAXPATH) < 0)
9      return -1;
10
11  begin_op();
12  if((ip = namei(old)) == 0){
13      end_op();
14      return -1;
15  }
16
17  ilock(ip);
18  if(ip->type == T_DIR){
19      iunlockput(ip);
20      end_op();
21      return -1;
22  }
23
24  ip->nlink++;
25  iupdate(ip);
26  iunlock(ip);
27
28  if((dp = nameiparent(new, name)) == 0)
29      goto bad;
30  ilock(dp);
31  if(dp->dev != ip->dev || dirlink(dp, name, ip->inum) < 0){
32      iunlockput(dp);
33      goto bad;
34  }
35  iunlockput(dp);
36  iput(ip);
37
38  end_op();
39
40  return 0;
41
42 bad:
43  ilock(ip);
44  ip->nlink--;
45  iupdate(ip);
46  iunlockput(ip);
47  end_op();
48  return -1;
49 }
```

## 修改 `open` 系统调用：

- 修改 `open` 系统调用，以处理路径指向符号链接的情况。如果文件不存在，`open` 应当失败。
- 如果 `open` 指定了 `O_NOFOLLOW` 标志，应当打开符号链接本身，而不是跟随符号链接。
- 如果符号链接指向的文件也是符号链接，必须递归地跟随直到达到非链接文件。如果符号链接形成了循环，应当返回错误代码。你可以通过设置一个递归深度的阈值（例如10）来近似处理这个问题。

```
1  uint64
2  sys_open(void)
3  {
4      char path[MAXPATH];
5      int fd, omode;
6      struct file *f;
7      struct inode *ip;
8      int n;
9
10     if((n = argstr(0, path, MAXPATH)) < 0 || argint(1, &omode) <
11        0)
12         return -1;
13     begin_op();
14
15     if(omode & O_CREATE){
16         ip = create(path, T_FILE, 0, 0);
17         if(ip == 0){
18             end_op();
19             return -1;
20         }
21     } else {
22         if((ip = namei(path)) == 0){
23             end_op();
24             return -1;
25         }
26         ilock(ip);
27         if(ip->type == T_DIR && omode != O_RDONLY){
28             iunlockput(ip);
29             end_op();
30             return -1;
```

```

31     }
32 }
33
34 if(ip->type == T_DEVICE && (ip->major < 0 || ip->major >=
NDEV)){
35     iunlockput(ip);
36     end_op();
37     return -1;
38 }
39
40 // handle the symlink - lab 9.2
41 if(ip->type == T_SYMLINK && (omode & O_NOFOLLOW) == 0) {
42     if((ip = follow_symlink(ip)) == 0) {
43         // 此处不用调用iunlockput()释放锁
44         // follow_symlinktest()返回失败时,锁在函数内已经被释放
45         end_op();
46         return -1;
47     }
48 }
49
50 if((f = filealloc()) == 0 || (fd = fdalloc(f)) < 0){
51     if(f)
52         fileclose(f);
53     iunlockput(ip);
54     end_op();
55     return -1;
56 }
57
58 if(ip->type == T_DEVICE){
59     f->type = FD_DEVICE;
60     f->major = ip->major;
61 } else {
62     f->type = FD_INODE;
63     f->off = 0;
64 }
65 f->ip = ip;
66 f->readable = !(omode & O_WRONLY);
67 f->writable = (omode & O_WRONLY) || (omode & O_RDWR);
68
69 if((omode & O_TRUNC) && ip->type == T_FILE){
70     itrunc(ip);
71 }

```

```

72
73     iunlock(ip);
74     end_op();
75
76     return fd;
77 }
78
79 // Generating symbolic links
80 uint64
81 sys_symlink(void) {
82     char target[MAXPATH], path[MAXPATH];
83     struct inode *ip;
84     int n;
85
86     if ((n = argstr(0, target, MAXPATH)) < 0
87         || argstr(1, path, MAXPATH) < 0) {
88         return -1;
89     }
90
91     begin_op();
92     // create the symlink's inode
93     if((ip = create(path, T_SYMLINK, 0, 0)) == 0) {
94         end_op();
95         return -1;
96     }
97     // write the target path to the inode
98     if(writei(ip, 0, (uint64)target, 0, n) != n) {
99         iunlockput(ip);
100         end_op();
101         return -1;
102     }
103
104     iunlockput(ip);
105     end_op();
106     return 0;
107 }

```

确保其他系统调用的正确性：

- 修改其他相关的系统调用（如 `link` 和 `unlink`），使它们操作符号链接本身，而不是跟随符号链接指向的文件。

```

1 // recursively follow the symlinks - lab9-2
2 // Caller must hold ip->lock
3 // and when function returned, it holds ip->lock of returned ip
4 static struct inode*
5 follow_symlink(struct inode* ip) {
6     uint inums[NSYMLINK];
7     int i, j;
8     char target[MAXPATH];
9
10    for(i = 0; i < NSYMLINK; ++i) {
11        inums[i] = ip->inum;
12        // read the target path from symlink file
13        if(readi(ip, 0, (uint64)target, 0, MAXPATH) <= 0) {
14            iunlockput(ip);
15            printf("open_symlink: open symlink failed\n");
16            return 0;
17        }
18        iunlockput(ip);
19
20        // get the inode of target path
21        if((ip = namei(target)) == 0) {
22            printf("open_symlink: path \"%s\" is not exist\n",
23                target);
24            return 0;
25        }
26        for(j = 0; j <= i; ++j) {
27            if(ip->inum == inums[j]) {
28                printf("open_symlink: links form a cycle\n");
29                return 0;
30            }
31            ilock(ip);
32            if(ip->type != T_SYMLINK) {
33                return ip;
34            }
35        }
36
37        iunlockput(ip);
38        printf("open_symlink: the depth of links reaches the
39            limit\n");
40        return 0;
41    }

```

实验得分