

Lecture 32: File system implementation in xv6

Mythili Vutukuru

IIT Bombay

<https://www.cse.iitb.ac.in/~mythili/os/>

Disk layout

- Disk in xv6 is formatted to contain superblock, log (for crash recovery), inode blocks (multiple inodes packed per block), bitmap (indicating which data blocks are free), actual data blocks
- Disk inode contains block numbers of direct blocks and one indirect block
- Directory is special file: data blocks contain directory entries, mapping file names of files in the directory to corresponding inode numbers
- Link count of inode = number of directory entries pointing to a file inode

```
4050 // On-disk file system format.
4051 // Both the kernel and user programs use this header file.
4052
4053
4054 #define ROOTINO 1 // root i-number
4055 #define BSIZE 512 // block size
4056
4057 // Disk layout:
4058 // [ boot block | super block | log | inode blocks |
4059 //                               free bit map | data blocks]
4060 //
4061 // mkfs computes the super block and builds an initial file system. The
4062 // super block describes the disk layout:
4063 struct superblock {
4064     uint size; // Size of file system image (blocks)
4065     uint nblocks; // Number of data blocks
4066     uint ninodes; // Number of inodes.
4067     uint nlog; // Number of log blocks
4068     uint logstart; // Block number of first log block
4069     uint inodestart; // Block number of first inode block
4070     uint bmapstart; // Block number of first free map block
4071 };
```

```
4077 // On-disk inode structure
4078 struct dinode {
4079     short type; // File type
4080     short major; // Major device number (T_DEV only)
4081     short minor; // Minor device number (T_DEV only)
4082     short nlink; // Number of links to inode in file system
4083     uint size; // Size of file (bytes)
4084     uint addrs[NDIRECT+1]; // Data block addresses
4085 };
```

```
4115 struct dirent {
4116     ushort inum;
4117     char name[DIRSIZ];
4118 };
```

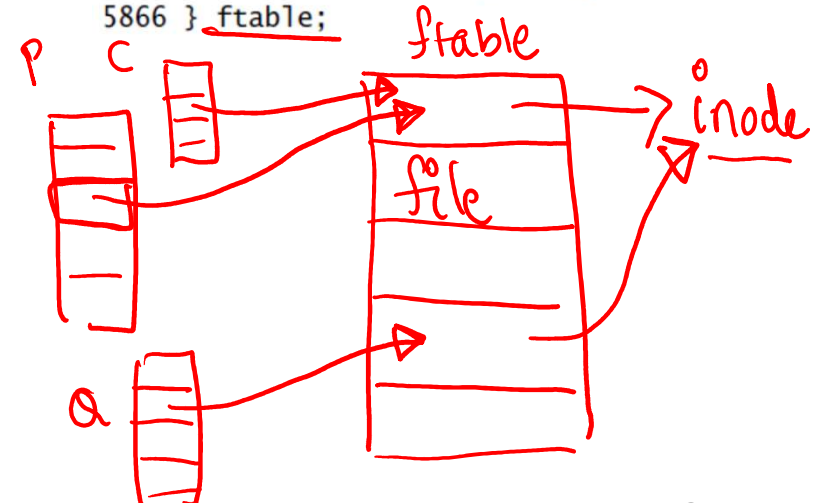
In-memory data structures (1)

- Every open file has a struct file associated with it
 - Pointer to inode or pipe structure
- All struct files stored in fixed size array called file table (ftable)
- File descriptor array of a process contains pointers to struct files in the file table
- Two processes P and Q open same file, will use two struct file entries in file table
 - Points to same inode
 - Read and write independently at different offsets
- P forks child C, both file descriptors will point to same struct file (ref is increased)
 - Offsets are shared
- Reference count of struct file is number of file descriptors that point to it

```
2348  struct file *ofile[NFILE];  // Open files
```

```
4150 struct file {  
4151     enum { FD_NONE, FD_PIPE, FD_INODE } type;  
4152     int ref; // reference count  
4153     char readable;  
4154     char writable;  
4155     struct pipe *pipe;  
4156     struct inode *ip;  
4157     uint off;  
4158 };
```

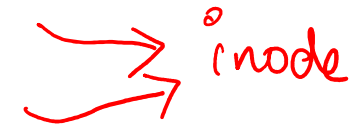
```
5862 struct devsw devsw[NDEV];  
5863 struct {  
5864     struct spinlock lock;  
5865     struct file file[NFILE];  
5866 } ftable;
```



In-memory data structures (2)

- Struct file points to in-memory inode structure of an open file (pipe structure for pipes)
- In-memory inode is almost copy of disk inode, stored in memory for open files
- All in-memory inodes stored in fixed size array called inode cache (icache)
- Reference count of in-memory inode is number of pointers from file table entries, current working directory of process etc.
 - Different from link count
 - A file is cleaned up on disk only when both ref count and link count are zero

```
4161 // in-memory copy of an inode
4162 struct inode {
4163     uint dev;           // Device number
4164     uint inum;          // Inode number
4165     int ref;            // Reference count
4166     struct sleeplock lock; // protects everything below here
4167     int valid;          // inode has been read from disk?
4168
4169     short type;         // copy of disk inode
4170     short major;
4171     short minor;
4172     short nlink;
4173     uint size;
4174     uint addrs[NDIRECT+1];
4175 };
```



```
5137 struct {
5138     struct spinlock lock;
5139     struct inode inode[NINODE];
5140 } icache;
```

Inode functions (1)



- Function ialloc() allocates free inode from disk by looking over disk inodes and finding a free one for a file
- Function iget() returns a reference counted pointer to in-memory inode in icache, to use in struct file etc
 - Non-exclusive pointer, information inside inode structure may not be up to date
 - Pointer released by iput()
- Function ilock() locks the inode for use by a process, and updates its information from disk if needed
 - Unlocked by iunlock()
- Function iupdate() propagates changes from in-memory inode to on-disk inode

ialloc
iget
ilock
iupdate

Inode functions (2)

- Inode has pointers to file datablocks
- Function `bmap` returns the address of n-th block of file
 - If direct block, read from inode
 - If indirect block, read indirect block first and then return block number from it
- Function can allocate data blocks too: if n-th data block of file not present, allocates new block on disk, writes it to inode, and returns address
- Functions readi/writei are used to read/write file data at given offset, call `bmap` to find corresponding data block

```
5400 // Inode content
5401 //
5402 // The content (data) associated with each inode is stored
5403 // in blocks on the disk. The first NDIRECT block numbers
5404 // are listed in ip->addrs[]. The next NINDIRECT blocks are
5405 // listed in block ip->addrs[NDIRECT].
5406
5407 // Return the disk block address of the nth block in inode ip.
5408 // If there is no such block, bmap allocates one.
5409 static uint
5410 bmap(struct inode *ip, uint bn)
5411 {
5412     uint addr, *a;
5413     struct buf *bp;
5414
5415     if(bn < NDIRECT){
5416         if((addr = ip->addrs[bn]) == 0)
5417             ip->addrs[bn] = addr = balloc(ip->dev);
5418         return addr;
5419     }
5420     bn -= NDIRECT;
5421
5422     if(bn < NINDIRECT){
5423         // Load indirect block, allocating if necessary.
5424         if((addr = ip->addrs[NDIRECT]) == 0)
5425             ip->addrs[NDIRECT] = addr = balloc(ip->dev);
5426         bp = bread(ip->dev, addr);
5427         a = (uint*)bp->data;
5428         if((addr = a[bn]) == 0){
5429             a[bn] = addr = balloc(ip->dev);
5430             log_write(bp);
5431         }
5432         brelse(bp);
5433         return addr;
5434     }
5435
5436     panic("bmap: out of range");
5437 }
```


Directory functions

```
4115 struct dirent {
4116     ushort inum;
4117     char name[DIRSIZ];
4118 };
```

- Directory lookup: read directory entries from the data blocks of directory. If file name matches, return pointer to inode from icache
- Linking a file to a directory: check file with the same name does not exist, and add the mapping from file name to inode number to directory

```
5608 // Look for a directory entry in a directory.
5609 // If found, set *poff to byte offset of entry.
5610 struct inode*
5611 dirlookup(struct inode *dp, char *name, uint *poff)
5612 {
5613     uint off, inum;
5614     struct dirent de;
5615
5616     if(dp->type != T_DIR)
5617         panic("dirlookup not DIR");
5618
5619     for(off = 0; off < dp->size; off += sizeof(de)){
5620         if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
5621             panic("dirlookup read");
5622         if(de.inum == 0)
5623             continue;
5624         if(namecmp(name, de.name) == 0){
5625             // entry matches path element
5626             if(poff)
5627                 *poff = off;
5628             inum = de.inum;
5629             return iget(dp->dev, inum);
5630         }
5631     }
5632
5633     return 0;
5634 }
```

```
5650 // Write a new directory entry (name, inum) into the directory dp.
5651 int
5652 dirlink(struct inode *dp, char *name, uint inum)
5653 {
5654     int off;
5655     struct dirent de;
5656     struct inode *ip;
5657
5658     // Check that name is not present.
5659     if((ip = dirlookup(dp, name, 0)) != 0){
5660         iput(ip);
5661         return -1;
5662     }
5663
5664     // Look for an empty dirent.
5665     for(off = 0; off < dp->size; off += sizeof(de)){
5666         if(readi(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
5667             panic("dirlink read");
5668         if(de.inum == 0)
5669             break;
5670     }
5671
5672     strncpy(de.name, name, DIRSIZ);
5673     de.inum = inum;
5674     if(writeti(dp, (char*)&de, off, sizeof(de)) != sizeof(de))
5675         panic("dirlink");
5676
5677     return 0;
5678 }
```

Creating a file (if it doesn't exist)

- Locate the inode of parent directory by walking the filepath from root (lookup root inode, find inode number of next element of pathname in inode data blocks, and repeat)
- Lookup filename in parent directory. If file already exists, return its inode
- If file doesn't exist, allocate a new inode for it, lock it, initialize it
- If new file is a directory, add entries for "." and ".."
- If new file is a regular file, link it to its parent directory

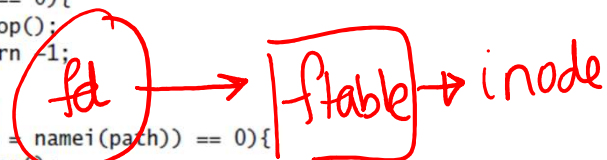
```
6356 static struct inode*
6357 create(char *path, short type, short major, short minor)
6358 {
6359     uint off;
6360     struct inode *ip, *dp;
6361     char name[DIRSIZ];
6362
6363     if((dp = nameiparent(path, name)) == 0)
6364         return 0;
6365     ilock(dp);
6366
6367     if((ip = dirlookup(dp, name, &off)) != 0){
6368         iunlockput(dp);
6369         ilock(ip);
6370         if(type == T_FILE && ip->type == T_FILE)
6371             return ip;
6372         iunlockput(ip);
6373         return 0;
6374     }
6375
6376     if((ip = ialloc(dp->dev, type)) == 0)
6377         panic("create: ialloc");
6378
6379     ilock(ip);
6380     ip->major = major;
6381     ip->minor = minor;
6382     ip->nlink = 1;
6383     iupdate(ip);
6384
6385     if(type == T_DIR){ // Create . and .. entries.
6386         dp->nlink++; // for "."
6387         iupdate(dp);
6388         // No ip->nlink++ for "..": avoid cyclic ref count.
6389         if(dirlink(ip, ".", ip->inum) < 0 || dirlink(ip, "..", dp->inum) < 0)
6390             panic("create dots");
6391     }
6392
6393     if(dirlink(dp, name, ip->inum) < 0)
6394         panic("create: dirlink");
6395
6396     iunlockput(dp);
6397
6398     return ip;
6399 }
```

/a/b/foolxt

System call: open

- Get arguments: filename, mode
- Create file (if specified) and get a pointer to its inode
- Allocate new struct file in ftable, and new file descriptor entry in struct proc of process pointing to the struct file in ftable
- Return index of new entry in file descriptor array of process
- Note the begin_op and end_op for transactions

```
6400 int
6401 sys_open(void)
6402 {
6403     char *path;
6404     int fd, omode;
6405     struct file *f;
6406     struct inode *ip;
6407
6408     if(argstr(0, &path) < 0 || argint(1, &omode) < 0)
6409         return -1;
6410
6411     begin_op();
6412
6413     if(omode & O_CREATE){
6414         ip = create(path, T_FILE, 0, 0);
6415         if(ip == 0){
6416             end_op();
6417             return -1;
6418         }
6419     } else {
6420         if((ip = namei(path)) == 0){
6421             end_op();
6422             return -1;
6423         }
6424         ilock(ip);
6425         if(ip->type == T_DIR && omode != O_RDONLY){
6426             iunlockput(ip);
6427             end_op();
6428             return -1;
6429         }
6430     }
6431
6432     if((f = filealloc()) == 0 || (fd = fdalloc(f)) < 0){
6433         if(f)
6434             fileclose(f);
6435         iunlockput(ip);
6436         end_op();
6437         return -1;
6438     }
6439     iunlock(ip);
6440     end_op();
6441
6442     f->type = FD_INODE;
6443     f->ip = ip;
6444     f->off = 0;
6445     f->readable = !(omode & O_WRONLY);
6446     f->writable = (omode & O_WRONLY) || (omode & O_RDWR);
6447     return fd;
6448 }
6449
```



```
graph LR
    fd((fd)) --> ftable[ftable]
    ftable --> inode[inode]
```

System call: link

ln old new

- Link an existing file from another directory with a new name (hard linking)
- Get pointer to file inode by walking the old filename
- Update link count in inode
- Get pointer to inode of new directory, and link old inode from parent directory in new name

```
6200 // Create the path new as a link to the same inode as old.
6201 int
6202 sys_link(void)
6203 {
6204     char name[DIRSIZ], *new, *old;
6205     struct inode *dp, *ip;
6206
6207     if(argstr(0, &old) < 0 || argstr(1, &new) < 0)
6208         return -1;
6209
6210     begin_op();
6211     if((ip = namei(old)) == 0){
6212         end_op();
6213         return -1;
6214     }
6215
6216     ilock(ip);
6217     if(ip->type == T_DIR){
6218         iunlockput(ip);
6219         end_op();
6220         return -1;
6221     }
6222
6223     ip->nlink++;
6224     iupdate(ip);
6225     iunlock(ip);
6226
6227     if((dp = nameiparent(new, name)) == 0)
6228         goto bad;
6229     ilock(dp);
6230     if(dp->dev != ip->dev || dirlink(dp, name, ip->inum) < 0){
6231         iunlockput(dp);
6232         goto bad;
6233     }
6234     iunlockput(dp);
6235     iput(ip);
6236
6237     end_op();
6238
6239     return 0;
6240
6241 bad:
6242     ilock(ip);
6243     ip->nlink--;
6244     iupdate(ip);
6245     iunlockput(ip);
6246     end_op();
6247     return -1;
6248 }
6249
```

new dir
*old ↙
inode*

System call: file read

- Other system calls follow same pattern
- For example, file read:
 - Get arguments (file descriptor number, buffer to read into, number of bytes to read)
 - Fetch inode pointer from struct file and perform read on inode (or pipe if file descriptor pointed to pipe)
 - Function readi uses the function “bmap” to get the block corresponding to n-th byte and reads from it
 - Offset in struct file updated

```
6131 int
6132 sys_read(void)
6133 {
6134     struct file *f;
6135     int n;
6136     char *p;
6137
6138     if(argfd(0, 0, &f) < 0 || argint(2, &n) < 0 || argptr(1, &p, n) < 0)
6139         return -1;
6140     return fileread(f, p, n);
6141 }
```

```
5963 // Read from file f.
5964 int
5965 fileread(struct file *f, char *addr, int n)
5966 {
5967     int r;
5968
5969     if(f->readable == 0)
5970         return -1;
5971     if(f->type == FD_PIPE)
5972         return piperead(f->pipe, addr, n);
5973     if(f->type == FD_INODE){
5974         ilock(f->ip);
5975         if((r = readi(f->ip, addr, f->off, n)) > 0)
5976             f->off += r;
5977         iunlock(f->ip);
5978         return r;
5979     }
5980     panic("fileread");
5981 }
```

Summary

- On disk: inodes, data blocks, free bitmap (and log)
- In-memory: file descriptor array (points to) struct file in file table array (points to) in-memory inode in inode cache
- Directory is a special file, where data blocks contain directory entries (filenames and corresponding inode numbers)
- System calls related to files extract arguments, perform various operations on in-memory and on-disk data structures
- Updates to disk happen via the buffer cache
 - Changes to all blocks in a system call are wrapped in a transaction and logged for atomicity