



Files and Directories

Referência principal

Ch.39 of *Operating Systems: Three Easy Pieces* by Remzi and Andrea Arpaci-Dusseau (pages.cs.wisc.edu/~remzi/OSTEP/)

Discutido em classe em 12 de novembro de 2018

How to Manage Persistent Devices?

How should the OS manage a persistent device?

What are the APIs for that?

What are the most important aspects of the implementation?

Persistent Storage

- A persistent device keeps data intact even if there is a power loss.
 - Hard disk drive
 - Solid-state storage device
- In a Unix-style system there are two key abstractions in the virtualization of persistent storage
 - File
 - Directory
- From a technical point of view there is no real difference between files and directories, although they play different roles in a system.

File

- A file is simply a linear array of bytes, which can be read or written.
- Every file or directory has an **inode**, which contains all the file's metadata (that is, all the administrative data needed to read a file is stored in its inode).
- Each file also has a low-level name, often called its **inode number**.
 - Usually, the user is not aware of this name.
- The file system is responsible for storing files persistently on disk and for making them available upon request.

Inodes

- An **inode** is an entry in an **inode table**, containing information (the metadata) about a regular file or directory.
 - Note that an inode does not store the name of the file.
- The **inode table** contains the **inodes** of all files in that file system.
 - The individual inodes in the inode table have a unique number (unique to that file system), the **inode number**.
- **An inode stores metadata about a file:**
 - Type: regular file, directory, pipe, etc.
 - Permissions: read, write, execute
 - Link count: number of hard links to it
 - User ID: owner
 - Group ID: group owner
 - Size of file: or major/minor number in case of some special files
 - Time stamps: access time, change time and (inode) change time
 - Attributes: 'immutable' for example
 - Access control list: permissions for special users/groups
 - Link to location of file
 - ...

Creating Files

- Creation is the most basic file operation.
- Use `open()` system call with `O_CREAT` flag.

```
int fd = open("foo", O_CREAT | O_WRONLY | O_TRUNC, S_IRUSR | S_IWUSR);
```

- `O_CREAT`: create file.
 - `O_WRONLY`: writing is the only operation allowed on that file.
 - `O_TRUNC`: make the file size zero (i.e. remove any existing content).
- File descriptor
 - A successful `open()` system call returns a file descriptor.
 - The descriptor is just an integer used to uniquely identify an open file of the process..

Open File Table

■ Open File Table

- A table indexed by file descriptors whose elements are pointers to file table entries.
- One unique open file table is provided by the operating system for each process.

■ Open File Table Entry

- An in-memory structure created when a process opens a file and whose components maintain file position and other indicators.

File I/O system calls

- In a Unix-style system, there are basically 5 types of file I/O system calls:
 - Create
 - Create a new empty file.
 - Open
 - Create a new file or open a file for reading, writing or both.
 - Close
 - To tell the OS you are done with a file descriptor and close the file which is pointed by it.
 - Read
 - Read a number of bytes from a file into a memory area.
 - Write
 - Write a number of bytes from a memory area into a file.

Standard file descriptors and calls

- When any process starts, file descriptors 0, 1 and 2 are created automatically.
 - By default the corresponding entries in the process' file table refer to a file named **/dev/tty**.
 - **/dev/tty** is an in-memory surrogate for the terminal, which is a combination of keyboard and video screen.

Standard file descriptors and calls

- System calls to standard files
 - **Read from stdin => read from fd 0**
 - Whenever we write any character from keyboard, it read from stdin through fd 0 and save to file named /dev/tty.
 - **Write to stdout => write to fd 1**
 - Whenever we see any output to the video screen, it's from the file named /dev/tty and written to stdout in screen through fd 1.
 - **Write to stderr => write to fd 2**
 - We see any error to the video screen, it is also from that file write to stderr in screen through fd 2.

Create a file

- Syntax in C language:

```
int creat(char *filename, mode_t mode);
```

- Parameters

- **filename**: name of file to be created
- **mode**: access permissions to new file

- Returns

- first unused file descriptor (generally 3 when first **creat** use in process, because 0, 1 and 2 are reserved)
- -1 when error

- How it works in OS

- Create new empty file on disk
- Create file table entry
- Set first unused file descriptor to point to file table entry
- Return file descriptor used, -1 upon failure

Open

- Syntax in C language

```
int open (const char* Path, int flags[, int mode ]);
```

- Parameters

- Path : path to file which you want to use
- use absolute path beginning with "/", when you are not working in same directory of file.
- Use relative path which is only file name with extension, when you are working in the same directory of file.

- Flags: how you would like to use the file

- `O_RDONLY`: read only
- `O_WRONLY`: write only
- `O_RDWR`: read and write
- `O_CREAT`: create file if it doesn't exist
- `O_EXCL`: prevent file creation if it already exists

- How it works in OS

- Find existing file on disk
- Create file table entry
- Set first unused file descriptor to point to file table entry
- Return file descriptor used, -1 upon failure

Close

- Syntax in C language

```
int close(int fd);
```

- Parameter

- `fd`: file descriptor

- Return

- 0 on success
- -1 on error

- How it works in the OS

- Destroy open file table entry referenced by element `fd`, as long as no other process is pointing to it
- Set element `fd` of file table to `NULL`

```
1.  #include <stdio.h>
2.  #include <stdlib.h>
3.  #include <unistd.h>
4.  #include <fcntl.h>

5.  int main(void) {
6.      int fd1 = open("foo.txt", O_CREAT | O_RDONLY, 0);
7.      if (fd1 < 0)
8.          { perror("c1"); exit(1); }
9.      printf("opened fd = %d\n", fd1);
10.
11.     if (close(fd1) < 0)
12.         { perror("c1"); exit(1); }
13.     printf("closed fd = %d\n", fd1);
14.
15.     int fd2 = open("bar.txt", O_CREAT | O_RDONLY, 0);
16.     printf("opened fd = %d\n", fd2);
17.     exit(0);
18. }
```

```
opened fd = 3
closed fd = 3
opened fd = 3
```

Read

- Syntax in C

```
size_t read(int fd, void* buf, size_t cnt);
```

- From the file indicated by `fd`, reads `cnt` bytes into the memory area indicated by `buf`.
 - If successful, it updates the access time for the file.
- Parameters
 - `fd`: file descriptor
 - `buf`: buffer to read data into
 - `cnt`: length of buffer
- Returns
 - Number of bytes read on success
 - 0 on reaching end of file
 - -1 on error
 - -1 on signal interrupt

Read

- Important points
 - `buf` must point to a valid memory location with length not smaller than the specified size to avoid overflow.
 - `fd` should be a valid file descriptor returned from `open()`
 - if `fd` is `NULL`, an error is generated.
 - `cnt` is the number of bytes to be read, while the return value is the actual number of bytes read.
 - Sometimes less than `cnt` bytes may be read.

Example of `read` system call

```
1. #include <stdio.h>
2. #include <stdlib.h>
3. #include <unistd.h>
4. #include <fcntl.h>

5. int main(void) {
6.     int fd, sz;
7.     char *c = (char *) calloc(100, sizeof(char));
8.
9.     fd = open("foo.txt", O_RDONLY);
10.    if (fd < 0)
11.        { perror("r1"); exit(1); }
12.
13.    sz = read(fd, c, 20);
14.    printf("System call read(%d, c, 20) returned that"
15.          " %d bytes were read.\n", fd, sz);
16.    c[sz] = '\0';
17.    printf("Those bytes were as follows: %s\n", c);
18. }
```

System call `read(3, c, 20)` returned that 10 bytes were read.
Those bytes were as follows: foobar 1 2

What is the output of the following program?

```
1. #include <stdio.h>
2. #include <stdlib.h>
3. #include <unistd.h>
4. #include <fcntl.h>
5.
6. int main() {
7.     char c1, c2;
8.     int fd1 = open("foo.txt", O_RDONLY, 0);
9.     int fd2 = open("foo.txt", O_RDONLY, 0);
10.    read(fd1, &c1, 1);
11.    read(fd2, &c2, 1);
12.    printf("c1 = %c\n", c1);
13.    printf("c2 = %c\n", c2);
14.    exit(0);
15. }
```

What is the output of the following program?

```
1. #include <stdio.h>
2. #include <stdlib.h>
3. #include <unistd.h>
4. #include <fcntl.h>
5.
6. int main() {
7.     char c1, c2;
8.     int fd1 = open("foo.txt", O_RDONLY, 0);
9.     int fd2 = open("foo.txt", O_RDONLY, 0);
10.    read(fd1, &c1, 1);
11.    read(fd2, &c2, 1);
12.    printf("c1 = %c\n", c1);
13.    printf("c2 = %c\n", c2);
14.    exit(0);
15. }
```

- The descriptors `fd1` and `fd2` each have their own file table entry, so each descriptor has its own file position for `foobar.txt`.
 - Thus, the read from `fd2` reads the first byte of `foobar.txt`, and the output is `c = f`, not `c = o`.

Write file

- Syntax in C

```
size_t write(int fd, void* buf, size_t cnt);
```

- Writes `cnt` bytes from `buf` to the file or socket associated with `fd`.
 - `cnt` should not be greater than `INT_MAX` (defined in the `limits.h` header file).
 - If `cnt` is zero, `write()` simply returns 0 without attempting any other action.
- Parameters
 - `fd`: file descriptor
 - `buf`: buffer to write data to
 - `cnt`: length of buffer
- Returns
 - Number of bytes written on success
 - 0 on reaching end of file
 - -1 on error or signal interrupt

Write file

- Important points
 - The file must be opened for write operations
 - `buf` needs to be at least as long as specified by `cnt` to prevent overflow.
 - `cnt` is the requested number of bytes to write, while the return value is the actual number of bytes written.
 - This happens when `fd` has a smaller than `cnt` number of bytes to write.
 - If `write()` is interrupted by a signal, the effect is one of the following:
 - If `write()` has not written any data yet, it returns -1 and sets `errno` to `EINTR`.
 - If `write()` has successfully written some data, it returns the number of bytes it wrote before it was interrupted.

Reading And Writing, But Not Sequentially

- An open file has a current offset which determines where the next read or write will start.
- The current offset can be updated
 - Implicitly
 - When a read or write of N bytes takes place, N is added to the current offset.
 - Explicitly
 - By calling `lseek()`.

Reading And Writing, But Not Sequentially

- Syntax in C

```
off_t lseek(int fildes, off_t offset, int whence);
```

- `fildes`: file descriptor
- `offset`: displacement from a location within the file
- `whence`: determine how offset is interpreted
- From the manual page...
 - If `whence` is `SEEK_SET`, the offset is set to `offset` bytes.
 - If `whence` is `SEEK_CUR`, the offset is set to its current location plus `offset` bytes.
 - If `whence` is `SEEK_END`, the offset is set to the size of the file plus `offset` bytes.

Writing immediately with `fsync()`

- The file system will buffer writes in memory for some time for performance reasons
- At that later point in time, the writes will actually be issued to the storage device.
 - Write seems to complete quickly but data can be lost, e.g. if the machine crashes.
- However, some applications require more than eventual guarantee.
 - E.g. DBMS requires forced writes to disk from time to time.

```
off_t fsync(int fd)
```

- File system forces all dirty (i.e., not yet written) data to disk for the file referred to by the file description.
- `fsync()` returns once all of these writes are complete.

Writing immediately with `fsync()`

- An example of `fsync()`

```
1.  int main(void) {  
2.      char buffer[12] = "hello world";  
3.      int size = strlen(buffer);  
4.      int fd = open("foo.txt", O_CREAT | O_WRONLY | O_TRUNC);  
5.      assert ( !(fd < 0) );  
6.      int rc = write(fd, buffer, size);  
7.      assert (rc == size);  
8.      rc = fsync(fd);  
9.      assert (rc == 0);  
10.     return 0;  
11. }
```

- In some cases, this may not be enough, and the code will also need to `fsync()` the directory that contains the file `foo.txt`.

Renaming Files

```
rename(char *old, char *new)
```

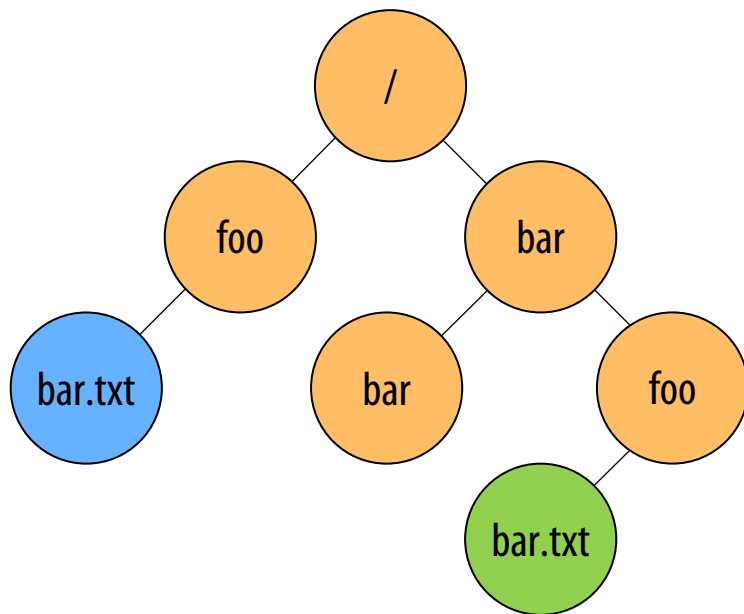
- Rename a file to a new name.
- It is implemented as an atomic call.

■ Example

```
1.  int main(void) {  
2.      char buffer[13] = "hello world\n";  
3.      int size = strlen(buffer);  
4.      int fd = open("rentest.tmp", O_WRONLY | O_CREAT | O_TRUNC);  
5.      write(fd, buffer, size); // write out new version of file  
6.      fsync(fd);  
7.      close(fd);  
8.      int renres = rename("rentest.tmp", "rentest.txt");  
9.      printf("rename returned %d\n", renres);  
10.     return 0;  
11. }
```

Directory

- A directory is like a file and also has a low-level name.
 - It contains a list of (user-readable name, low-level name) pairs.
 - Each entry in a directory refers to either a file or another directory.



- By placing directories within other directories, users can build an arbitrary directory tree (or directory hierarchy), under which files and other directories can be stored.
- The directory hierarchy starts at the root directory (named / in Unix-based systems) and uses a separator to create a chain of directory names until the desired file is reached.
 - The resulting compound name is called the file's absolute pathname.
 - The last part of the file name, named extension, usually indicates the type of the file.

Creating Directories

```
1. ~$ mkdir tstdir
2. ~$ cd tstdir
3. tstdir$ mkdir foo
4. tstdir$ ls
5. foo
6. tstdir$ cd foo
7. foo$ ls
8. foo$ ls -a
9. . .
10. foo$ ls -ail
11. total 0
12. 105312409 drwxr-xr-x  2 arthur.catto  staff  64 Nov 12 17:06 .
13. 105312376 drwxr-xr-x  3 arthur.catto  staff  96 Nov 12 17:06 ..
14. foo$
```

A directory is created empty.

An empty directory has two entries:
self and parent

Each associated to a different inode

Reading Directories

```
1.  int main(void) {
2.      DIR *dp = opendir(".");    // open current directory
3.      assert(dp != NULL);
4.      struct dirent *d;
5.      while ((d = readdir(dp)) != NULL) { // read one directory entry
6.          // print out the name and inode number of each file
7.          printf("%d %s\n", (int) d->d_ino, d->d_name);
8.      }
9.      closedir(dp);              // close current directory
10.     return 0;
11. }
```

while the information available within `struct dirent` is

```
1.  struct dirent {
2.      char d_name[256];          /* filename */
3.      ino_t d_ino;                /* inode number */
4.      off_t d_off;                /* offset to the next direct */
5.      unsigned short d_reclen;    /* length of this record */
6.      unsigned char d_type;       /* type of file */
7.  }
```

Hard Links

- The **identity** of a file is its inode number, not its name.
- A **hard link** is a name that references an inode.
 - If *file1* has a hard link named *file2*, then both refer to same inode.
- Creating a hard link for a file means adding a new name to an inode.
 - This can be done by the `link` function or the `ln` command.

```
int link(const char *path1, const char *path2);
```

```
1. ~/tstdir$ cd foo
2. ~/tstdir/foo$ ls
3. ~/tstdir/foo$ echo "hello world" > file1
4. ~/tstdir/foo$ cat file1
5. hello world
6. ~/tstdir/foo$ ln file1 file2
7. ~/tstdir/foo$ cat file2
8. hello world
9. ~/tstdir/foo$ rm file1
10. ~/tstdir/foo$ cat file2
11. hello world
12. ~/tstdir/foo$
```

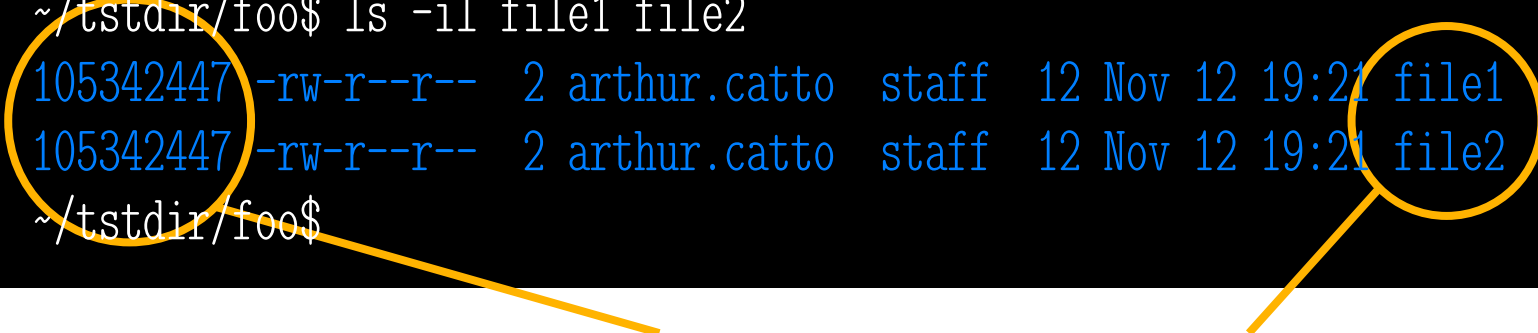
Hard Links

- This is how `link(file1, file2)` works
 - Create another name in the directory.
 - Refer it to the same inode number of the original file.
 - The file is not copied in any way.
 - Then, we now just have two human names (`file1` and `file2`) that both refer to the same file.

Hard Links

- This is the result of `link()`

```
1. ~/tstdir/foo$ ls -il file1 file2
2. 105342447 -rw-r--r--  2 arthur.catto  staff  12 Nov 12 19:21 file1
3. 105342447 -rw-r--r--  2 arthur.catto  staff  12 Nov 12 19:21 file2
4. ~/tstdir/foo$
```



- Two files have the same inode number, but two human names.
 - There is no difference between `file1` and `file2`.
 - Both just link to the underlying metadata about the file.

Symbolic Links

- A **symbolic link** is a separate file whose contents point to the linked-to file.
 - To create a symbolic link, use the `ln` command with the option `-s`.
- Symbolic links are more flexible than hard links.
 - We cannot create a hard link to a directory.
 - We cannot create a hard link to a file in another partition.
 - Can you tell why?
 - Because inode numbers are only unique within a file system.

Symbolic links

```
1. ~/tstdir/foo$ ls
2. file1 file2
3. ~/tstdir/foo$ ln -s file1 file3
4. ~/tstdir/foo$ ls
5. file1 file2 file3
6. ~/tstdir/foo$ ls -l
7. total 16
8. -rw-r--r--  2 arthur.catto  staff  12 Nov 12 19:21 file1
9. -rw-r--r--  2 arthur.catto  staff  12 Nov 12 19:21 file2
10. lrwxr-xr-x  1 arthur.catto  staff   5 Nov 12 19:48 file3 -> file1
11. ~/tstdir/foo$ ls -il
12. total 16
13. 105342447 -rw-r--r--  2 arthur.catto  staff  12 Nov 12 19:21 file1
14. 105342447 -rw-r--r--  2 arthur.catto  staff  12 Nov 12 19:21 file2
15. 105350240 lrwxr-xr-x  1 arthur.catto  staff   5 Nov 12 19:48 file3 -> file1
```

Symbolic links

```
1. ~/tstdir/foo$ cat file1
2. hello world
3. ~/tstdir/foo$ cat file2
4. hello world
5. ~/tstdir/foo$ cat file3
6. hello world
7. ~/tstdir/foo$ ls -ails
8. total 16
9. 105312409 0 drwxr-xr-x  5 arthur.catto  staff  160 Nov 12 19:48 .
10. 105312376 0 drwxr-xr-x  3 arthur.catto  staff   96 Nov 12 17:06 ..
11. 105342447 8 -rw-r--r--  2 arthur.catto  staff   12 Nov 12 19:21 file1
12. 105342447 8 -rw-r--r--  2 arthur.catto  staff   12 Nov 12 19:21 file2
13. 105350240 0 lrwxr-xr-x  1 arthur.catto  staff    5 Nov 12 19:48 file3 -> file1
14. ~/tstdir/foo$
```

Symbolic links

```
1. ~/tstdir/foo$ echo "hello world" > longerfilenameA
2. ~/tstdir/foo$ ln -s longerfilenameA fileB
3. ~/tstdir/foo$ ls -ails
4. total 24
5. 105312409 0 drwxr-xr-x  9 arthur.catto  staff  288 Nov 12 20:12 .
6. 105312376 0 drwxr-xr-x  3 arthur.catto  staff   96 Nov 12 17:06 ..
7. 105354762 8 -rw-r--r--   1 arthur.catto  staff  12 Nov 12 20:12 longerfilenameA
8. 105354350 0 lrwxr-xr-x   1 arthur.catto  staff   5 Nov 12 20:10 areallyveryverylongnewname -> file1
9. 105354151 0 lrwxr-xr-x   1 arthur.catto  staff   5 Nov 12 20:08 averylongnewname -> file1
10. 105342447 8 -rw-r--r--   2 arthur.catto  staff  12 Nov 12 19:21 file1
11. 105342447 8 -rw-r--r--   2 arthur.catto  staff  12 Nov 12 19:21 file2
12. 105350240 0 lrwxr-xr-x   1 arthur.catto  staff   5 Nov 12 19:48 file3 -> file1
13. 105354911 0 lrwxr-xr-x   1 arthur.catto  staff  16 Nov 12 20:12 fileB -> longerfilenameA
14. ~/tstdir/foo$
```

Removing files

- To remove a file, we call `unlink()`.

```
1. ~/tstdir/foo$ rm fileB
2. ~/tstdir/foo$ rm file2
3. ~/tstdir/foo$ ls -ails
4. total 16
5. 105312409 0 drwxr-xr-x 7 arthur.catto staff 224 Nov 12 20:18 .
6. 105312376 0 drwxr-xr-x 3 arthur.catto staff 96 Nov 12 17:06 ..
7. 105354762 8 -rw-r--r-- 1 arthur.catto staff 12 Nov 12 20:12 alongerfilenameA
8. 105354350 0 lrwxr-xr-x 1 arthur.catto staff 5 Nov 12 20:10 areallyveryverylongnewname ->
   file1
9. 105354151 0 lrwxr-xr-x 1 arthur.catto staff 5 Nov 12 20:08 averylongnewname -> file1
10. 105342447 8 -rw-r--r-- 1 arthur.catto staff 12 Nov 12 19:21 file1
11. 105350240 0 lrwxr-xr-x 1 arthur.catto staff 5 Nov 12 19:48 file3 -> file1
12. ~/tstdir/foo$
```

- Reference count

- Tracks how many different file names have been linked to this inode.
- When `unlink()` is called, the reference count is decremented.
- If the reference count reaches zero, the file system frees the inode and related data blocks, actually deleting the file.

Deleting Directories

- To delete a directory we use `rmdir()`.
 - It requires the directory to be empty.
 - If we try to remove a non-empty directory, the call will fail.