

## Reference:

Advanced Programming in the UNIX Environment, Third Edition, W. Richard Stevens and Stephen A. Rago, Addison-Wesley Professional Computing Series, 2013. Chapter 3

*Computer Systems: A Programmer's Perspective*, Randal E. Bryant and David R. O'Hallaron, Prentice Hall, 3<sup>rd</sup> edition, 2016, Chapter 10

Acknowledgment: Several examples are taken from CSAPP book cited above

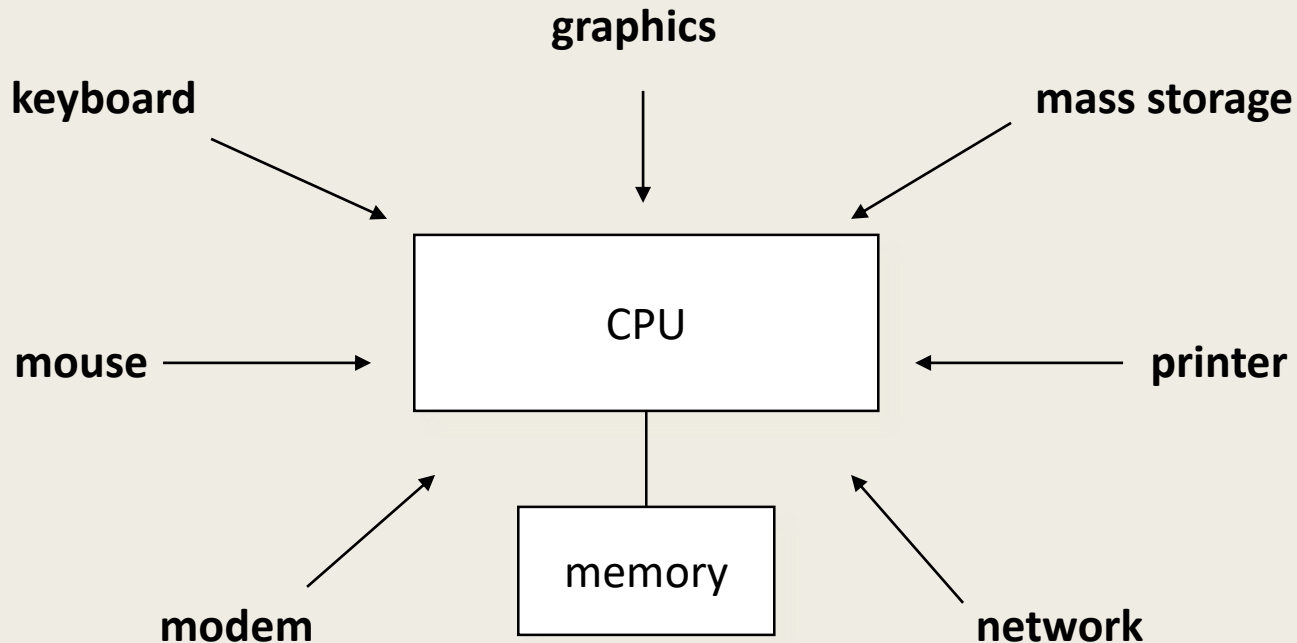
# UNIX I/O

Tanzir Ahmed  
CSCE 313 Spring 2020

# What is I/O

- I/O is the process of copying data between main memory and external devices (disk drives, terminals, networks)
- Language run-time systems provide higher-level facilities for performing I/O (e.g. `printf`, `scanf`)
  - *On Unix systems, these higher-level I/O functions eventually call **System-Level Unix I/O functions** provided by the Kernel (e.g., ***printf()*** internally calls ***write()***)*
  - *The same applies to Windows and other OSs*
    - `printf()` calls native windows function `WriteFile()`

# Files are not always “Files”: I/O Devices



# Unix File Types

- **Regular files**

- *File containing user/app data (binary, text, whatever)*

- **Directory files**

- *A file that contains the names and locations of other files*

- **Character special and block special files**

- *Terminals (character special) and disks (block special)*

- **FIFO (named pipe)**

- *A file type used for inter-process communication*

- **Socket**

- *A file type used for network communication between processes*

# Unix I/O

## ■ Key Features

- *Elegant mapping of files to devices allows kernel a simple interface called Unix I/O*
- *Important idea: All input and output is handled in a consistent and uniform way*

## ■ Basic Unix I/O operations (system calls):

- *Opening and closing files*

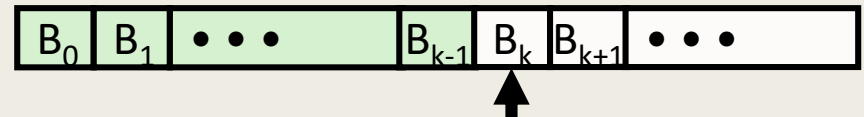
- **open()** and **close()**

- *Reading and writing a file*

- **read()** and **write()**

- **Current file position**

- Kernel maintains a file position (initially 0) for each open file (aka **file pointer**)
- Indicates next (byte) offset into file to read or write
- Reading and writing automatically **advance** the file pointer
- **lseek()** can be used to change file position **manually** at programmer's wish
  - *Instead of the default which is tied to read and write operations*



Current file position =  $k$

# Opening Files

- Opening a file informs the kernel that you are getting ready to access that file

```
int fd;    /* file descriptor */
if ((fd = open("/etc/hosts", O_RDONLY)) < 0) {
    perror("open");
    exit(1);
}
```

- Returns an identifying integer *file descriptor*
  - *fd == -1* indicates that an error occurred
- Kernel keeps track of all information about the open file. E.g., current file position, permissions etc.
- Each process created by a Unix shell begins with three open files associated with the terminal:
  - *0: standard input*
  - *1: standard output*
  - *2: standard error*
- On success, open returns the smallest available fd value
  - E.g., in the above *fd == 3* (0-2 are occupied)

# Closing Files

- Closing a file informs the kernel that you are finished accessing that file

```
int fd;      /* file descriptor */
int retval; /* return value */

if ((retval = close(fd)) < 0) {
    perror("close");
    exit(1);
}
```

- The kernel responds by
  - *freeing some data structures*
  - *restoring the descriptor to the pool of available descriptors*
- When a process terminates for any reason, the kernel performs **close** on all open files

# Example

- What is the output of the following program?

```
int main () {  
    int fd1 = open("foo.txt", O_RDONLY);  
    close(fd1);  
    int fd2 = open("baz.txt", O_RDONLY);  
    printf("fd2=%d\n", fd2);  
}
```

- Unix processes begin life with open descriptors assigned to stdin (fd=0), stdout (fd=1), and stderr (fd=2).
- The open function always returns the lowest unopened descriptor so the output will be **"fd2=3"**



# Reading Files

- Reading a file copies bytes from the current file position to memory, and then updates file position

```
char buf[512];
int fd;          /* file descriptor */
int nbytes;      /* number of bytes read */

/* Open file fd ... */
/* Then read up to 512 bytes from file fd */
if ((nbytes = read(fd, buf, sizeof(buf))) < 0) {
    perror("read");
    exit(1);
}
```

- Returns number of bytes read from file `fd` into `buf`
  - Return type **`ssize_t`** is signed integer
  - **`nbytes < 0`** indicates that an error occurred
  - **Short counts** (**`nbytes < sizeof(buf)`**) are possible and are not errors (eg. EOF when the file does not have `nbytes` of data starting from the file pointer, reading text from terminal where user gave `< nbytes` etc.)

# Reading Files - Example

- Output of the following program reading a file containing “foobar”

```
char buf [2];  
int fd;          /* file descriptor */  
int nbytes;      /* number of bytes read */  
  
/* Open file fd ... */  
nbytes = read(fd, buf, sizeof(buf)); // buf = ?  
nbytes = read(fd, buf, sizeof(buf)); // buf = ?
```

**Answer: “fo” and “ob”**  
**Why not “fo” and “foob”??**

# Writing Files

- Writing a file **writes** bytes from memory to the current file position, and then **updates** current file position

```
char buf[512];
int fd;          /* file descriptor */
int nbytes;      /* number of bytes written*/

/* Open the file fd ... */
/* Then write up to 512 bytes from buf to file fd */
if ((nbytes = write(fd, buf, sizeof(buf))) < 0) {
    perror("write");
    exit(1);
}
```

- Returns number of bytes written from `buf` to file `fd`
  - ***nbytes < 0*** indicates that an error occurred
  - As with reads, short counts are possible and are not errors!

# Writing Files - Example

- What would be the state of the file after executing the following program?

```
char buf[] = {'a', 'b'}; // what if char* buf = "ab"
int fd;           /* file descriptor */
int nbytes;       /* number of bytes written*/

/* Open the file fd ... */
/* Then write up to 2 bytes from buf to file fd */
write(fd, buf, sizeof(buf)); //file = ?
write(fd, buf, sizeof(buf)); //file = ?
```

abab

# File Representation – 3 Different Tables

## ■ Descriptor Table

- *Each process has its own separate descriptor table (only accessible directly by the kernel)*
- *Entries are indexed by the process's open file descriptors*
- *Each entry points to a **File Table** entry*

## ■ File Table

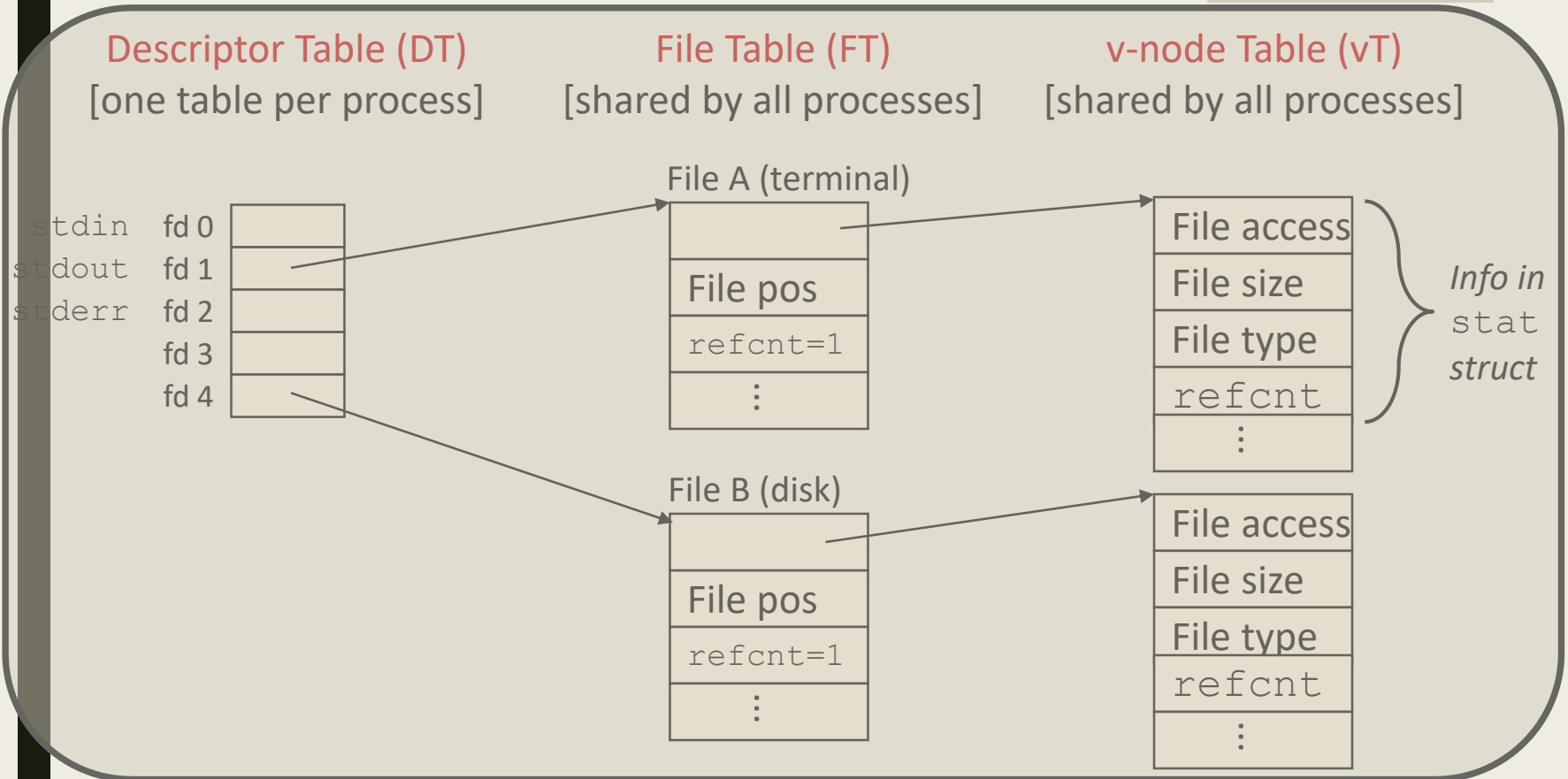
- *Shared by all processes.*
- *Each entry consists of*
  - **File Cursor**
  - **Reference count:** # of descriptor entries from all processes that point to it (why do we need this??)
  - ptr to the v-node table entry, along with status flags

## ■ v-node Table

- *Each entry (1 per file) contains information in the stat structure*

# How the Unix Kernel Represents Open Files

Two descriptors referencing two distinct open disk files. Descriptor 1 (stdout) points to terminal, and descriptor 4 points to open file in kernel memory



# File Sharing

- Two distinct DTs sharing the same disk file through two distinct open FT entries
  - E.g., calling **open** twice with the same **filename** argument*
  - Every **open** creates a new DT and a new FT entry, also possibly a vT entry*

Descriptor table (DT)

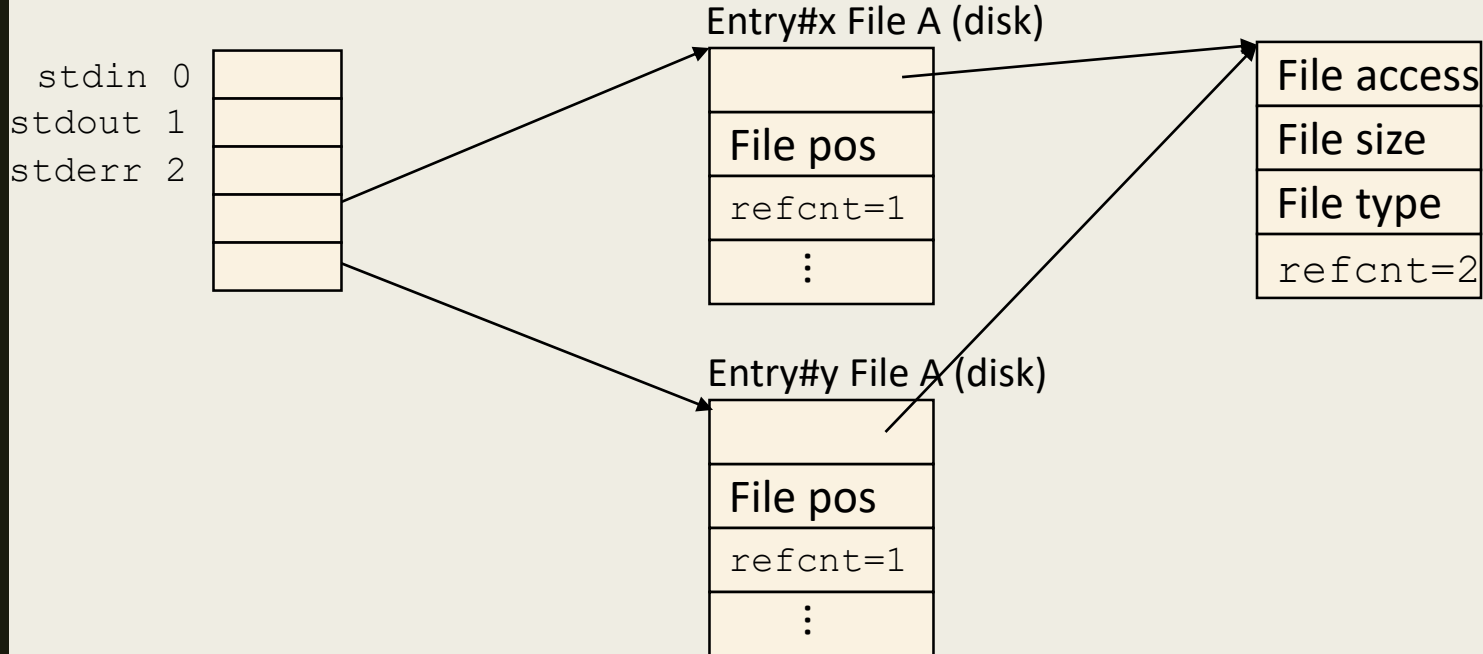
[one table per process]

Open File Table (FT)

[shared by all processes]

v-node Table (vT)

[shared by all processes]



# Example

Suppose the disk file foobar.txt consists of the six ASCII characters “foobar”. Then what is the output of the following program:

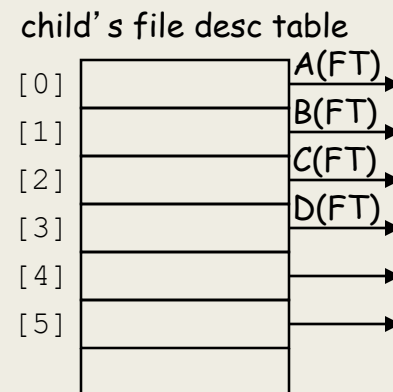
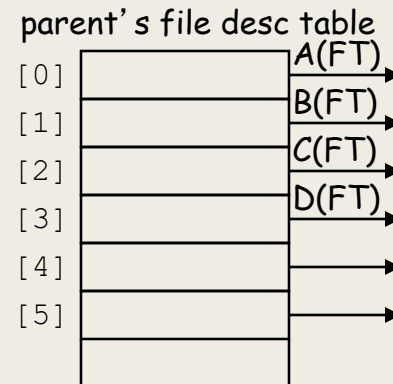
```
int main () {
    int fd1, fd2;
    char c;
    fd1 = open("foobar.txt", O_RDONLY, 0);
    fd2 = open("foobar.txt", O_RDONLY, 0);
    read(fd1, &c, 1);
    printf("c=%c\n", c);
    read(fd2, &c, 1);
    printf("c=%c\n", c);
}
```

*Answer: The descriptors fd1 and fd2 each have their own open file table entry so each descriptor has its own file position for foobar.txt. Thus the read from fd1 reads the first byte of foobar.txt and the output is “c=f”. Same is true for fd2 and the output “c=f”*



# File Descriptors and `fork()`

With `fork()`, child inherits content of parent's address space, including file descriptor table



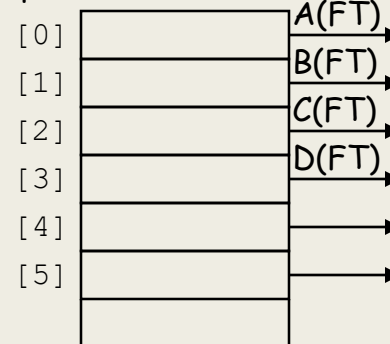
system file table (FT)

|               |
|---------------|
|               |
| A             |
| B             |
|               |
| C             |
| D ("myf.txt") |
|               |
|               |

# File Descriptors and `fork()`

```
int main(void) {  
    char c;  
    int myfd =  
open("myf.txt", O_RDONLY);  
    fork();  
    read(myfd, &c, 1);  
    printf("Got %c\n", c);  
}
```

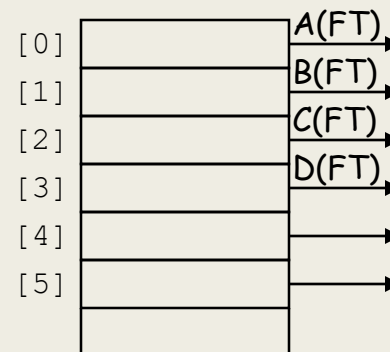
parent's file desc table



system file table (FT)

|               |
|---------------|
|               |
| A             |
| B             |
|               |
| C             |
| D ("myf.txt") |
|               |
|               |

child's file desc table

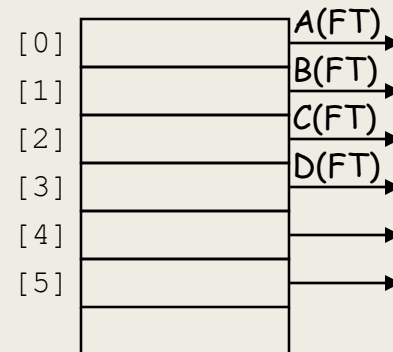


# File Descriptors and `fork()`

## (III)

```
int main(void) {  
    char c;  
    fork();  
    int myfd =  
    open("myf.txt", O_RDONLY);  
    read(myfd, &c, 1);  
    printf("Got %c\n", c);  
}
```

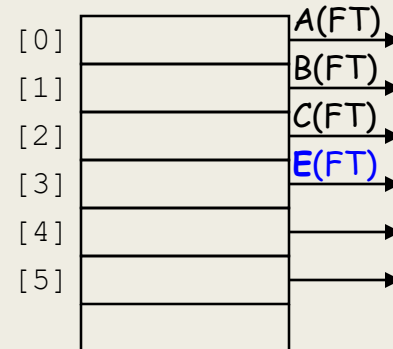
parent's file desc table



system file table (FT)

|               |
|---------------|
|               |
| A             |
| B             |
|               |
| C             |
| D ("myf.txt") |
|               |
|               |
| E ("myf.txt") |
|               |
|               |

child's file desc table



# Example

Suppose the disk file foobar.txt consists of the six ASCII characters “foobar”. Then what is the output of the following program:

```
int main () {  
    int fd;  
    char c;  
    fd = open("foobar.txt", O_RDONLY, 0);  
    if (fork() == 0) {  
        read(fd, &c, 1);  
        exit(0);  
    }  
    wait(NULL)  
    read(fd, &c, 1);  
    printf("c=%c\n", c);  
    exit(0);  
}
```

*Answer: The child inherits the parent's descriptor table and all processes share the same file table. Thus the descriptor fd in both the parent and child points to the same open file table entry. When the child reads the first byte of the file, the file position increments by 1. Thus the parent reads the second byte and output is “c=o”*

# I/O Redirection

- Question: How does a shell implement I/O redirection?

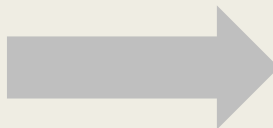
***unix> ls > foo.txt***

- Answer: By calling the `dup2 (oldfd, newfd)` function
  - Copies (per-process) descriptor table entry ***oldfd*** to entry ***newfd***

Descriptor table

*before* `dup2 (4, 1)`

|      |   |
|------|---|
| fd 0 |   |
| fd 1 | a |
| fd 2 |   |
| fd 3 |   |
| fd 4 | b |



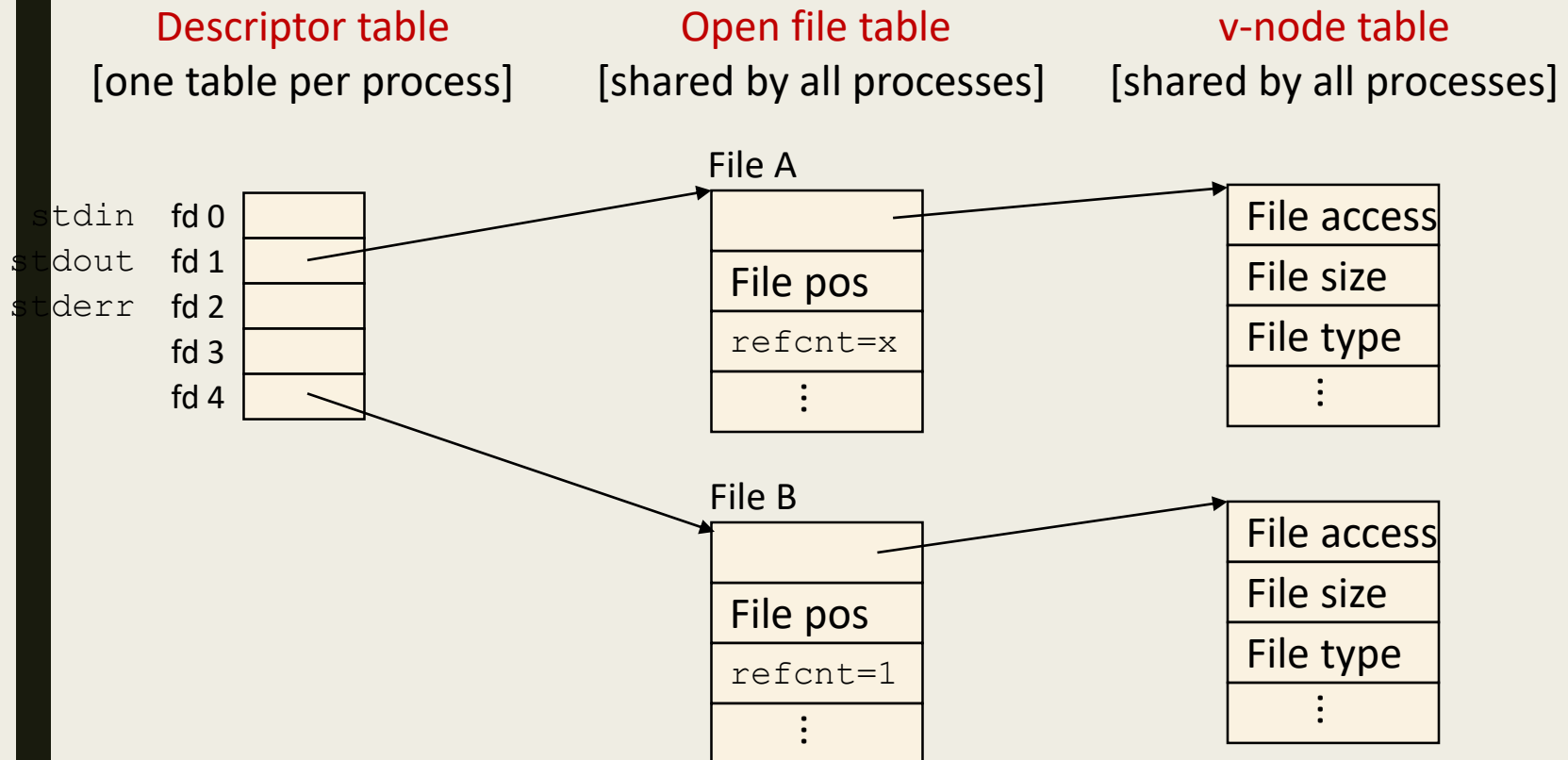
Descriptor table

*after* `dup2 (4, 1)`

|      |   |
|------|---|
| fd 0 |   |
| fd 1 | b |
| fd 2 |   |
| fd 3 |   |
| fd 4 | b |

# I/O Redirection Example

Step #1: open a that will be the target of redirection



# /0 Redirection Example (cont.)

Step #2: call `dup2 (4, 1)` to redirect stdout to 4

- cause `fd=1` (stdout) to refer to disk file pointed at by `fd=4`

**Descriptor table**

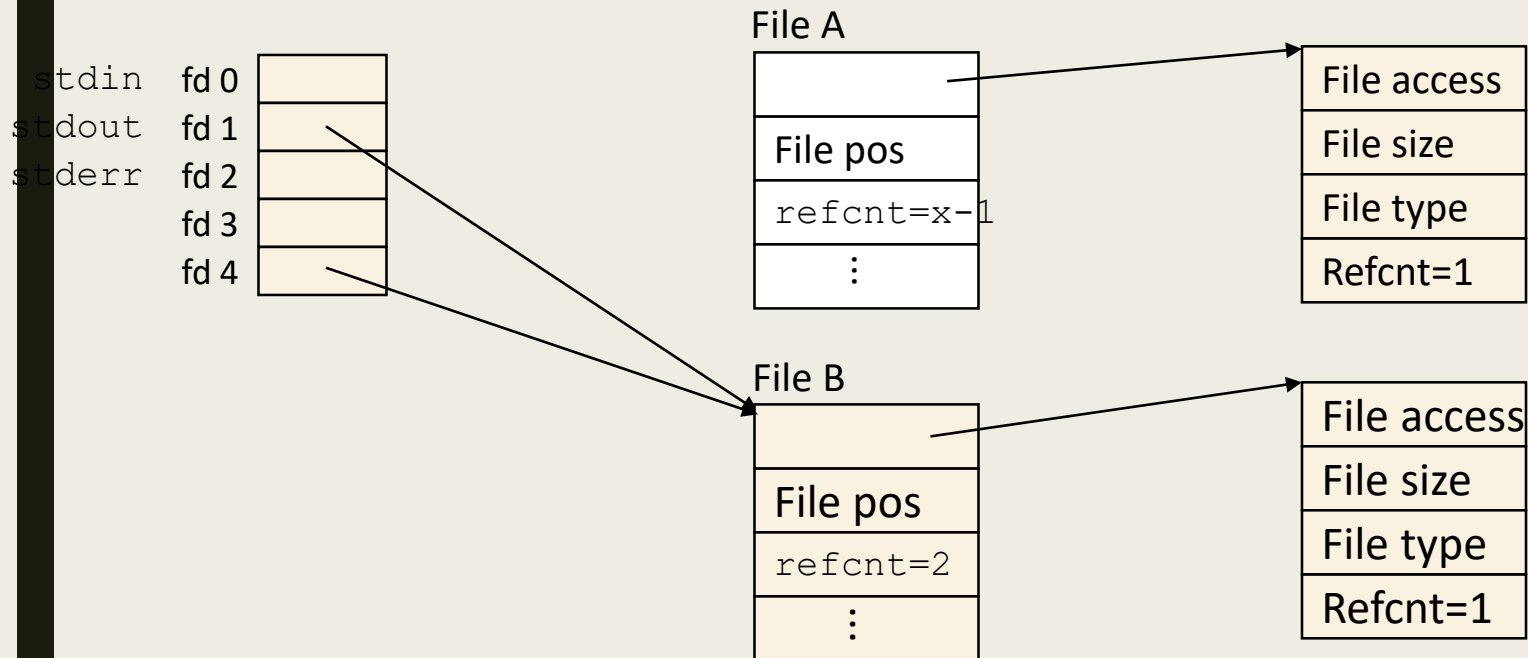
[one table per process]

**Open file table**

[shared by all processes]

**v-node table**

[shared by all processes]



# Example

Assuming that the disk file foobar.txt consists of six ASCII characters “foobar” what is the output?

```
int main(){
    int fd1, fd2;
    char c;
    char * fname = "foobar.txt";
    fd1 = open(fname, O_RDONLY, 0);
    fd2 = open(fname, O_RDONLY, 0);
    read(fd2, &c, 1);
    dup2(fd2, fd1);
    read(fd1, &c, 1);
    printf("c = %c\n", c);
    return 0;
}
```

- Because we are redirecting fd1 to fd2, the output is c=o



# Practice: Fun with File Descriptors (1)

```
int main(int argc, char *argv[])
{
    int fd1, fd2, fd3;
    char c1, c2, c3;
    char *fname = argv[1];
    fd1 = open(fname, O_RDONLY, 0);
    fd2 = open(fname, O_RDONLY, 0);
    fd3 = open(fname, O_RDONLY, 0);
    dup2(fd2, fd3);
    read(fd1, &c1, 1);
    read(fd2, &c2, 1);
    read(fd3, &c3, 1);
    printf("c1 = %c, c2 = %c, c3 = %c\n", c1, c2, c3);
    return 0;
}
```

- What would this program print for file containing “abcde”?

# Practice: Fun with File Descriptors (2)

```
int main(int argc, char *argv[])
{
    int fd1, fd2, fd3;
    char *fname = argv[1];
    fd1 = open(fname, O_CREAT|O_TRUNC|O_RDWR,
S_IRUSR|S_IWUSR);
    write(fd1, "pqrs", 4);
    fd2 = open(fname, O_APPEND|O_WRONLY, 0);
    write(fd2, "jklmn", 5);
    fd3 = dup(fd1); /* Allocates descriptor */
    write(fd3, "wxyz", 4);
    write(fd2, "ef", 2);
    // playing with file position
    lseek (fd1, 0, SEEK_SET);
    write (fd1, "ab", 2);
    return 0;
}
```

# Standard I/O Functions

- The C standard library (`libc.so`) contains a collection of higher-level *standard I/O* functions
  - *Documented in Appendix B of K&R.*
- Examples of standard I/O functions:
  - *Opening and closing files (**fopen** and **fclose**)*
  - *Reading and writing bytes (**fread** and **fwrite**)*
  - *Reading and writing text lines (**fgets** and **fputs**)*
  - *Formatted reading and writing (**fscanf** and **fprintf**)*

# Standard I/O Streams

- Standard I/O models open files as *streams*
  - *Abstraction for a file descriptor and a buffer in memory.*
- C programs begin life with three open streams (defined in `stdio.h`)
  - ***stdin*** (*standard input*)
  - ***stdout*** (*standard output*)
  - ***stderr*** (*standard error*)

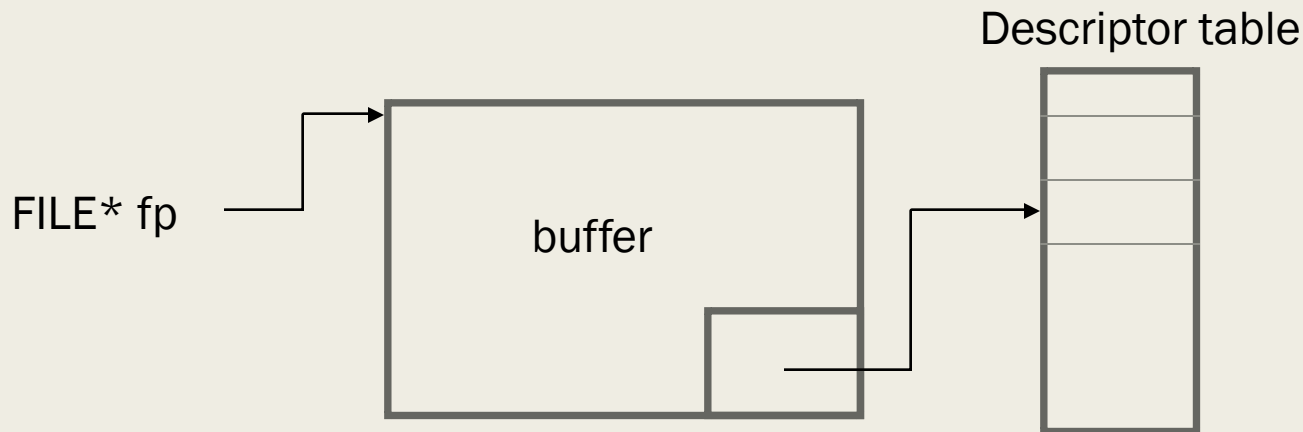
```
#include <stdio.h>
extern FILE *stdin; /* standard input (descriptor 0) */
extern FILE *stdout; /* standard output (descriptor 1) */
extern FILE *stderr; /* standard error (descriptor 2) */

int main() {
    fprintf(stdout, "Hello, world\n");
}
```

# Standard I/O – An Extra Layer

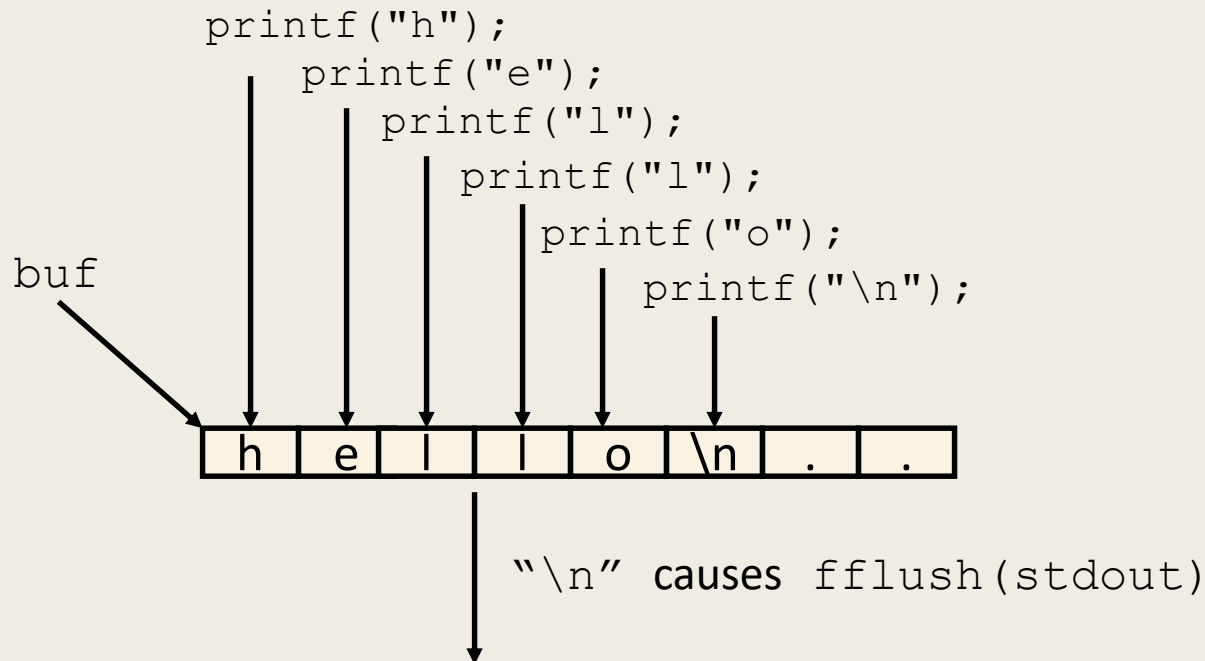
## ■ Motivations

- Portability: Code is compiled to Windows-native or Linux-native code, not tied to a particular architecture
- Buffering: For efficiency, considering device type
  - **fprintf** to **stdout** is flushed on ‘\n’ character
  - **fprintf** to disk files is flushed after reaching file buffer size



# Buffering in Standard I/O

- Standard I/O functions use buffered I/O



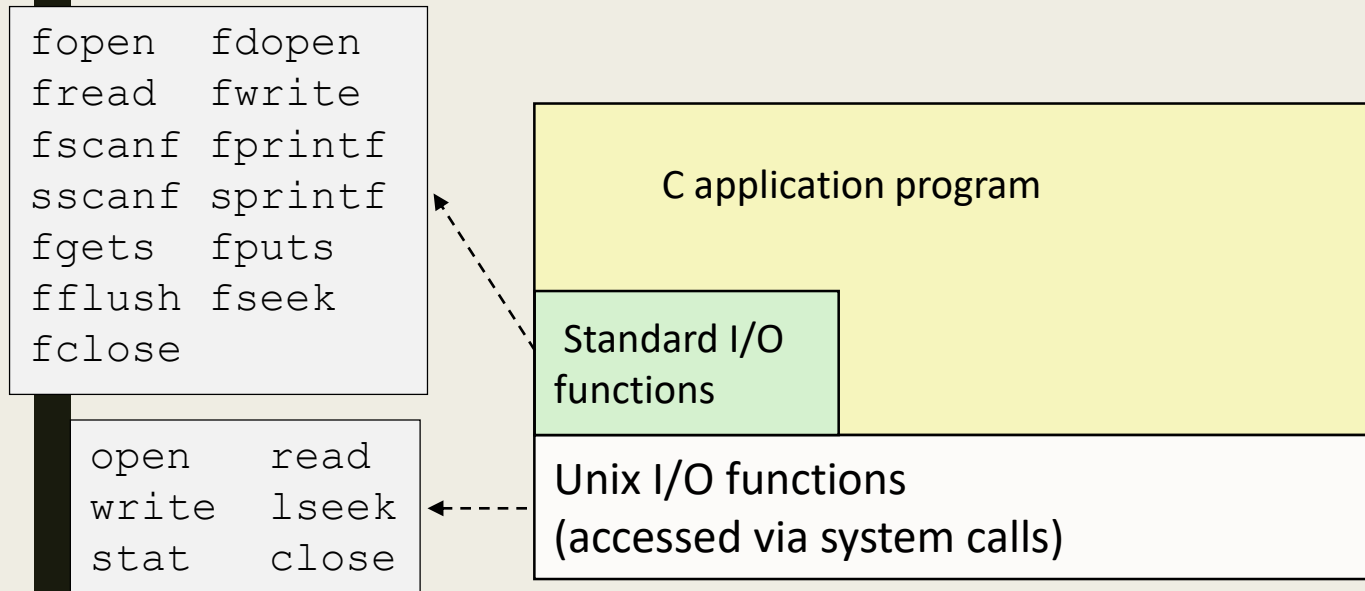
■

```
write(1, buf, 6);
```

- Buffer flushed to output fd on `"\n"` or `fflush()` call

# Unix I/O vs. Standard I/O

- Standard I/O is implemented using low-level Unix I/O



- Which ones should you use in your programs?

# Pros and Cons of Unix I/O

## ■ Pros

- *It is the most general and lowest overhead form of I/O*
  - All other I/O packages are implemented using Unix I/O functions on a Unix system
- *It provides functions for accessing file metadata*

## ■ Cons

- *Efficient data access requires some form of buffering, which is:*
  - Device dependent (e.g., a whole track of a disk can be read at once)
  - tricky and error prone
- *Both of these issues are addressed by standard I/O packages*



# Pros and Cons of Standard I/O

## ■ Pros:

- *Portable code – same code works on Windows and Unix*
- *Buffering increases efficiency by decreasing the number of **read** and **write** system calls*
- *Less burden on the programmer*

## ■ Cons:

- *Provides no function for accessing file metadata*
- *Standard I/O is not appropriate for input and output on network sockets*
  - *There are poorly documented restrictions on streams that interact badly with restrictions on sockets (CS:APP2e, Sec 10.9)*

# Choosing I/O Functions

- General rule: use the highest-level I/O functions you can
  - *Many C programmers are able to do all of their work using the standard I/O functions*
- When to use standard I/O
  - *When working with disk or terminal files*
- When to use raw Unix I/O
  - *Inside signal handlers, because Unix I/O is async-signal-safe*
  - *In rare cases when you need absolute highest performance*
    - You do your own buffering depending on the application nature and knowledge about the underlying hardware
    - On a 24 hard-drive (i.e., 48 TB cap) RAID system, we needed to exploit windows-native I/O to get nearly 3 GB/s read/write speed
    - Standard I/O would only give us < 1GBps read/write speed

# For Further Information

- The Unix bible:
  - *W. Richard Stevens & Stephen A. Rago, **Advanced Programming in the Unix Environment**, 3<sup>rd</sup> Edition, Addison Wesley, 2013*
- *Computer Systems: A Programmer's Perspective*, Randal E. Bryant and David R. O'Hallaron,
  - *Prentice Hall, 3<sup>rd</sup> edition, 2016, Chapter 10*