

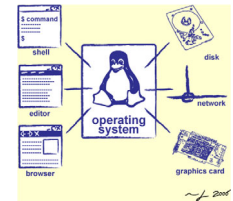
CS162
Operating Systems and
Systems Programming
Lecture 4

Abstractions 2: Files and I/O
A quick programmer's viewpoint

September 9th, 2020
Prof. John Kubiawicz
<http://cs162.eecs.Berkeley.edu>

Goals for Today: The File Abstraction

- Finish discussion of process management
- High-Level File I/O: Streams
- Low-Level File I/O: File Descriptors
- *How* and *Why* of High-Level File I/O
- Process State for File Descriptors
- Common Pitfalls with OS Abstractions



9/9/20

Kubiawicz CS162 © UCB Fall 2020

Lec 4.2

Recall: Synchronization between threads

- **Mutual Exclusion:** Ensuring only one thread does a particular thing at a time (one thread *excludes* the others)
- **Critical Section:** Code that exactly one thread can execute at once
 - Result of mutual exclusion
- **Lock:** An object only one thread can hold at a time
 - **Provides** mutual exclusion
- Offers two **atomic** operations:
 - Lock.Acquire() – wait until lock is free; then grab
 - Lock.Release() – Unlock, wake up waiters
- Need other tools for “cooperation”
 - e.g., semaphores

Semaphores: A quick look

- Semaphores are a kind of *generalized lock*
 - First defined by Dijkstra in late 60s
 - Main synchronization primitive used in original UNIX (& Pintos)
 - Definition: a Semaphore has a non-negative integer value and supports the following two operations:
 - **P()** or **down()**: atomic operation that waits for semaphore to become positive, then decrements it by 1
 - **V()** or **up()**: an atomic operation that increments the semaphore by 1, waking up a waiting P, if any
- P()** stands for “*proberen*” (to test) and **V()** stands for “*verhogen*” (to increment) in Dutch

9/9/20

Kubiawicz CS162 © UCB Fall 2020

Lec 4.3

9/9/20

Kubiawicz CS162 © UCB Fall 2020

Lec 4.4

Two Semaphore Patterns

- **Mutual Exclusion:** (like lock)
 - Called a "binary semaphore" or "mutex"
 - initial value of semaphore = 1;
 - semaphore.down();
 - // Critical section goes here
 - semaphore.up();
- **Signaling** other threads, e.g. **ThreadJoin**

Initial value of semaphore = 0

```
ThreadJoin {
    semaphore.down();
}
ThreadFinish {
    semaphore.up();
}
```

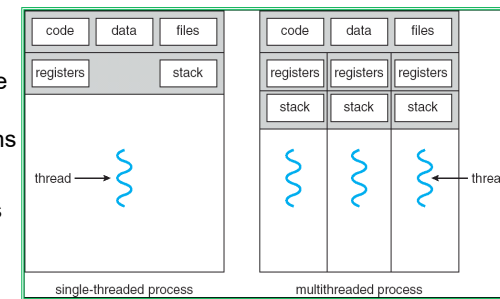
9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.5

Recall: Processes

- Definition: execution environment with restricted rights
 - One or more threads executing in a single address space
 - Owns file descriptors, network connections
- Instance of a running program
 - When you run an executable, it runs in its own process
 - Application: one or more processes working together
- Protected from each other; OS protected from them
- **In modern OSes, anything that runs outside of the kernel runs in a process**



9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.6

Recall: Creating Processes

- **pid_t fork()** – copy the current process
 - New process has different pid
 - New process contains a single thread
- Return value from **fork()**: pid (like an integer)
 - When > 0:
 - » Running in (original) **Parent** process
 - » return value is **pid** of new child
 - When = 0:
 - » Running in new **Child** process
 - When < 0:
 - » Error! Must handle somehow
 - » Running in original process
- **State of original process duplicated in both Parent and Child!**
 - Address Space (Memory), File Descriptors (covered later), etc...

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.7

Recall: fork_race.c

```
int i;
pid_t cpid = fork();
if (cpid > 0) {
    for (i = 0; i < 10; i++) {
        printf("Parent: %d\n", i);
        // sleep(1);
    }
} else if (cpid == 0) {
    for (i = 0; i > -10; i--) {
        printf("Child: %d\n", i);
        // sleep(1);
    }
} else { /* ERROR! */ }
```

Parent Process
Runs HERE!

Child Process
Runs HERE!

- What does this print?
- Would adding the calls to `sleep()` matter?

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.8

Recall: Start new Program with exec

```
...
cpid = fork();
if (cpid > 0) {
    /* Parent Process */
    tcpid = wait(&status);
} else if (cpid == 0) {
    /* Child Process */
    char *args[] = {"ls", "-l", NULL};
    execv("/bin/ls", args);

    /* execv doesn't return when it works.
       So, if we got here, it failed! */

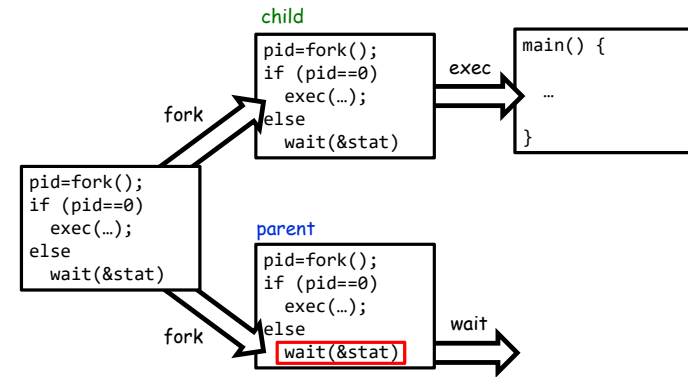
    perror("execv");
    exit(1);
}
...
```

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.9

Starting New Program (for instance in Shell)



9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.10

Finishing up: Process Management API

- `exit` – terminate a process
- `fork` – copy the current process
- `exec` – change the *program* being run by the current process
- `wait` – wait for a process to finish
- `kill` – send a *signal* (interrupt-like notification) to another process
- `sigaction` – set handlers for signals

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.11

fork2.c – parent waits for child to finish

```
int status;
pid_t tcpid;
...
cpid = fork();
if (cpid > 0) {
    /* Parent Process */
    mypid = getpid();
    printf("[%d] parent of [%d]\n", mypid, cpid);
    tcpid = wait(&status);
    printf("[%d] bye %d(%d)\n", mypid, tcpid, status);
} else if (cpid == 0) {
    /* Child Process */
    mypid = getpid();
    printf("[%d] child\n", mypid);
    exit(42);
}
...
```

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.12

Finishing up: Process Management API

- `exit` – terminate a process
- `fork` – copy the current process
- `exec` – change the *program* being run by the current process
- `wait` – wait for a process to finish
- `kill` – send a *signal* (interrupt-like notification) to another process
- `sigaction` – set handlers for signals

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.13

inf_loop.c

```
#include <stdlib.h>
#include <stdio.h>
#include <sys/types.h>
#include <unistd.h>
#include <signal.h>

void signal_callback_handler(int signum) {
    printf("Caught signal!\n");
    exit(1);
}

int main() {
    struct sigaction sa;
    sa.sa_flags = 0;
    sigemptyset(&sa.sa_mask);
    sa.sa_handler = signal_callback_handler;
    sigaction(SIGINT, &sa, NULL);
    while (1) {}
}
```

Q: What would happen if the process receives a SIGINT signal, but does not register a signal handler?

A: The process dies!

For each signal, there is a default handler defined by the system

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.14

Common POSIX Signals

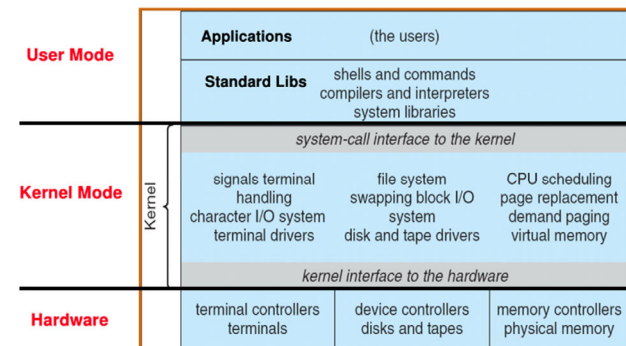
- SIGINT – control-C
- SIGTERM – default for `kill` shell command
- SIGSTP – control-Z (default action: stop process)
- SIGKILL, SIGSTOP – terminate/stop process
 - Can't be changed with `sigaction`
 - Why?

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.15

Recall: UNIX System Structure

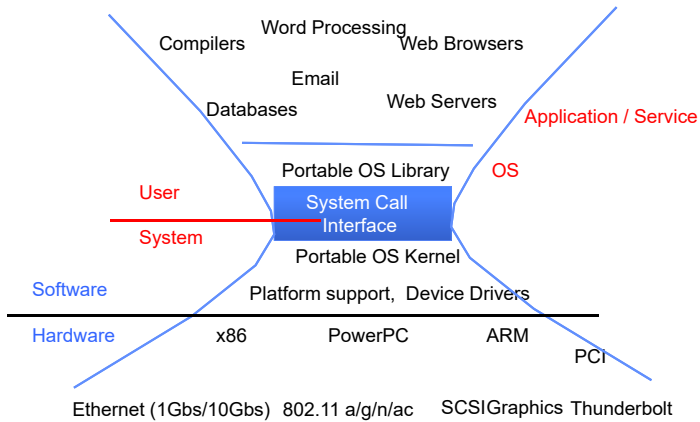


9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.16

Recall: System Calls (“Syscalls”)

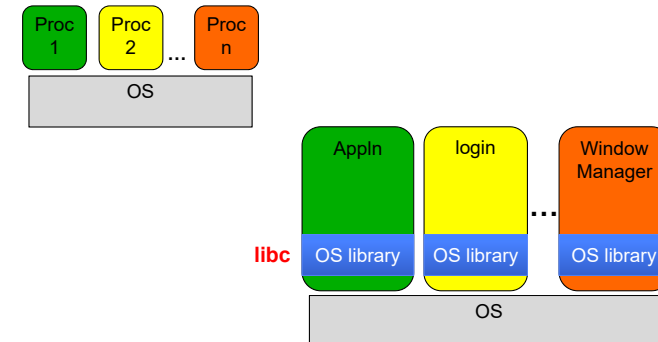


9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.17

Recall: OS Library Issues Syscalls



9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.18

Administrivia: Game On!

- **Project 0 due today!**
 - To be done on your own – like a homework!
- Slip days: I'd bank these and not spend them right away!
 - No credit when late and run out of slip days
- Group assignment in process (preferences were due Monday night)
 - Plan on attending your permanent discussion section this Friday
 - Remember to turn on your camera in Zoom
 - Discussion section attendance is mandatory
- **Midterm 1: October 1st, 5-7PM (Three weeks from tomorrow!)**
 - We understand that this partially conflicts with CS170, but those of you in CS170 can start that exam after 7PM (according to CS170 staff)
 - Video Proctored, No curve, Use of computer to answer questions
 - More details as we get closer to exam
- **Start Planning on how your group will collaborate on projects!**
 - Virtual Coffee Hours with your group (with camera)
 - Regular Brainstorming meetings?
 - Try to meet multiple times a week



9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.19

What does *pthread* stand for?

- pthread library: **POSIX** thread library
- **POSIX: Portable Operating System Interface (for uniX?)**
 - Interface for application programmers (mostly)
 - Defines the term “Unix,” derived from AT&T Unix
 - Created to bring order to many Unix-derived OSes, so applications are portable
 - » Partially available on non-Unix OSes, like Windows
 - Requires standard system call interface

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.20

Unix/POSIX Idea: Everything is a “File”

- Identical interface for:
 - Files on disk
 - Devices (terminals, printers, etc.)
 - Regular files on disk
 - Networking (sockets)
 - Local interprocess communication (pipes, sockets)
- Based on the system calls **open()**, **read()**, **write()**, and **close()**
- Additional: **ioctl()** for custom configuration that doesn't quite fit
- Note that the “Everything is a File” idea was a radical idea when proposed
 - Dennis Ritchie and Ken Thompson described this idea in their seminal paper on UNIX called “The UNIX Time-Sharing System” from 1974
 - I posted this on the resources page if you are curious

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.21

The File System Abstraction

- File
 - Named collection of data in a file system
 - POSIX File data: sequence of bytes
 - » Could be text, binary, serialized objects, ...
 - File Metadata: information about the file
 - » Size, Modification Time, Owner, Security info, Access control
- Directory
 - “Folder” containing files & directories
 - Hierarchical (graphical) naming
 - » Path through the directory graph
 - » Uniquely identifies a file or directory
 - /home/ff/cs162/public_html/fa14/index.html
 - Links and Volumes (later)

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.22

Connecting Processes, File Systems, and Users

- **Every process has a *current working directory (CWD)***
 - Can be set with system call:


```
int chdir(const char *path); //change CWD
```
- Absolute paths ignore CWD
 - /home/oski/cs162
- Relative paths are relative to CWD
 - index.html, ./index.html
 - » Refers to index.html in current working directory
 - ../index.html
 - » Refers to index.html in parent of current working directory
 - ~/index.html, ~cs162/index.html
 - » Refers to index.html in the home directory

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.23

I/O and Storage Layers



9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.24

Today: The File Abstraction

- High-Level File I/O: Streams
- Low-Level File I/O: File Descriptors
- *How* and *Why* of High-Level File I/O
- Process State for File Descriptors
- Common Pitfalls with OS Abstractions

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.25

C High-Level File API – Streams

- Operates on “streams” – unformatted sequences of bytes (with text or binary data), with a position: 

```
#include <stdio.h>
FILE *fopen( const char *filename, const char *mode );
int fclose( FILE *fp );
```

Mode	Text	Binary	Descriptions
r		rb	Open existing file for reading
w		wb	Open for writing; created if does not exist
a		ab	Open for appending; created if does not exist
r+		rb+	Open existing file for reading & writing.
w+		wb+	Open for reading & writing; truncated to zero if exists, create otherwise
a+		ab+	Open for reading & writing. Created if does not exist. Read from beginning, write as append

- Open stream represented by **pointer** to a **FILE** data structure
 - Error reported by returning a NULL pointer

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.25

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.26

C API Standard Streams – `stdio.h`

- Three predefined streams are opened implicitly when the program is executed.
 - FILE `*stdin` – normal source of input, can be redirected
 - FILE `*stdout` – normal source of output, can too
 - FILE `*stderr` – diagnostics and errors
- STDIN / STDOUT enable composition in Unix
- All can be redirected
 - `cat hello.txt | grep "World!"`
 - `cat`'s `stdout` goes to `grep`'s `stdin`

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.27

C High-Level File API

```
// character oriented
int fputc( int c, FILE *fp );           // rtn c or EOF on err
int fputs( const char *s, FILE *fp );  // rtn > 0 or EOF

int fgetc( FILE *fp );
char *fgets( char *buf, int n, FILE *fp );

// block oriented
size_t fread( void *ptr, size_t size_of_elements,
              size_t number_of_elements, FILE *a_file );
size_t fwrite( const void *ptr, size_t size_of_elements,
              size_t number_of_elements, FILE *a_file );

// formatted
int fprintf( FILE *restrict stream, const char *restrict format, ... );
int fscanf( FILE *restrict stream, const char *restrict format, ... );
```

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.27

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.28

C Streams: Char-by-Char I/O

```
int main(void) {
    FILE* input = fopen("input.txt", "r");
    FILE* output = fopen("output.txt", "w");
    int c;

    c = fgetc(input);
    while (c != EOF) {
        fputc(output, c);
        c = fgetc(input);
    }
    fclose(input);
    fclose(output);
}
```

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.29

C High-Level File API

```
// character oriented
int fputc( int c, FILE *fp );      // rtn c or EOF on err
int fputs( const char *s, FILE *fp ); // rtn > 0 or EOF

int fgetc( FILE * fp );
char *fgets( char *buf, int n, FILE *fp );

// block oriented
size_t fread(void *ptr, size_t size_of_elements,
             size_t number_of_elements, FILE *a_file);
size_t fwrite(const void *ptr, size_t size_of_elements,
             size_t number_of_elements, FILE *a_file);

// formatted
int fprintf(FILE *restrict stream, const char *restrict format, ...);
int fscanf(FILE *restrict stream, const char *restrict format, ... );
```

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.30

C Streams: Block-by-Block I/O

```
#define BUFFER_SIZE 1024
int main(void) {
    FILE* input = fopen("input.txt", "r");
    FILE* output = fopen("output.txt", "w");
    char buffer[BUFFER_SIZE];
    size_t length;
    length = fread(buffer, BUFFER_SIZE, sizeof(char), input);
    while (length > 0) {
        fwrite(buffer, length, sizeof(char), output);
        length = fread(buffer, BUFFER_SIZE, sizeof(char), input);
    }
    fclose(input);
    fclose(output);
}
```

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.31

Aside: System Programming

- Systems programmers should always be paranoid!
 - Otherwise you get intermittently buggy code
- We should really be writing things like:

```
FILE* input = fopen("input.txt", "r");
if (input == NULL) {
    // Prints our string and error msg.
    perror("Failed to open input file")
}
```
- Be thorough about checking return values!
 - Want failures to be systematically caught and dealt with
- I may be a bit loose with error checking for examples in class (to keep short)
 - Do as I say, not as I show in class!

9/9/20

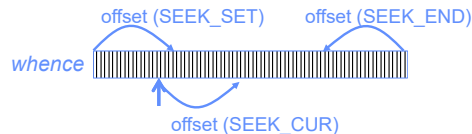
Kubiatowicz CS162 © UCB Fall 2020

Lec 4.32

C High-Level File API: Positioning The Pointer

```
int fseek(FILE *stream, long int offset, int whence);
long int ftell (FILE *stream)
void rewind (FILE *stream)
```

- For `fseek()`, the offset is interpreted based on the `whence` argument (constants in `stdio.h`):
 - `SEEK_SET`: Then offset interpreted from beginning (position 0)
 - `SEEK_END`: Then offset interpreted backwards from end of file
 - `SEEK_CUR`: Then offset interpreted from current position



- Overall preserves high-level abstraction of a uniform stream of objects

Today: The File Abstraction

- High-Level File I/O: Streams
- Low-Level File I/O: File Descriptors
- *How* and *Why* of High-Level File I/O
- Process State for File Descriptors
- Common Pitfalls with OS Abstractions [if time]

Key Unix I/O Design Concepts

- Uniformity – everything is a file
 - file operations, device I/O, and interprocess communication through `open`, `read/write`, `close`
 - Allows simple composition of programs
 - » `find | grep | wc ...`
- Open before use
 - Provides opportunity for access control and arbitration
 - Sets up the underlying machinery, i.e., data structures
- Byte-oriented
 - Even if blocks are transferred, addressing is in bytes
- Kernel buffered reads
 - Streaming and block devices looks the same, read blocks yielding processor to other task
- Kernel buffered writes
 - Completion of out-going transfer decoupled from the application, allowing it to continue
- Explicit close

Low-Level File I/O: The RAW system-call interface

```
#include <fcntl.h>
#include <unistd.h>
#include <sys/types.h>

int open (const char *filename, int flags [, mode_t mode])
int creat (const char *filename, mode_t mode)
int close (int filedes)
```

Bit vector of:

- Access modes (Rd, Wr, ...)
- Open Flags (Create, ...)
- Operating modes (Appends, ...)

Bit vector of Permission Bits:

- User|Group|Other X R|W|X

- Integer return from `open()` is a *file descriptor*
 - *Error indicated by return < 0: the global `errno` variable set with error (see man pages)*
- Operations on *file descriptors*:
 - Open system call created an *open file description* entry in system-wide table of open files
 - *Open file description* object in the kernel represents an instance of an open file
 - *Why give user an integer instead of a pointer to the file description in kernel?*

C Low-Level (pre-opened) Standard Descriptors

```
#include <unistd.h>
STDIN_FILENO - macro has value 0
STDOUT_FILENO - macro has value 1
STDERR_FILENO - macro has value 2

// Get file descriptor inside FILE *
int fileno (FILE *stream)

// Make FILE * from descriptor
FILE * fdopen (int filedes, const char *opentype)
```

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.37

Low-Level File API

- Read data from open file using file descriptor:

```
ssize_t read (int filedes, void *buffer, size_t maxsize)
```

 - Reads up to maxsize bytes – **might actually read less!**
 - returns bytes read, 0 => EOF, -1 => error
- Write data to open file using file descriptor

```
ssize_t write (int filedes, const void *buffer, size_t size)
```

 - returns number of bytes written
- Reposition file offset within kernel (this is independent of any position held by high-level FILE descriptor for this file!)

```
off_t lseek (int filedes, off_t offset, int whence)
```

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.38

Example: lowio.c

```
int main() {
    char buf[1000];
    int    fd = open("lowio.c", O_RDONLY, S_IRUSR | S_IWUSR);
    ssize_t rd = read(fd, buf, sizeof(buf));
    int    err = close(fd);
    ssize_t wr = write(STDOUT_FILENO, buf, rd);
}
```

- How many bytes does this program read?

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.39

POSIX I/O: Design Patterns

- **Open before use**
 - Access control check, setup happens here
- **Byte-oriented**
 - Least common denominator
 - OS responsible for hiding the fact that real devices may not work this way (e.g. hard drive stores data in blocks)
- **Explicit close**

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.40

POSIX I/O: Kernel Buffering

- **Reads are buffered inside kernel**
 - Part of making everything byte-oriented
 - Process is **blocked** while waiting for device
 - Let other processes run while gathering result
- **Writes are buffered inside kernel**
 - Complete in background (more later on)
 - Return to user when data is “handed off” to kernel
- This buffering is part of global buffer management and caching for block devices (such as disks)
 - Items typically cached in quanta of disk block sizes
 - We will have many interesting things to say about this buffering when we dive into the kernel

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.41

Low-Level I/O: Other Operations

- Operations specific to terminals, devices, networking, ...
 - e.g., `ioctl`
- Duplicating descriptors
 - `int dup2(int old, int new);`
 - `int dup(int old);`
- Pipes – channel
 - `int pipe(int pipefd[2]);`
 - Writes to `pipefd[1]` can be read from `pipefd[0]`
- File Locking
- Memory-Mapping Files
- Asynchronous I/O

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.42

Today: The File Abstraction

- High-Level File I/O: Streams
- Low-Level File I/O: File Descriptors
- **How and Why of High-Level File I/O**
- Process State for File Descriptors
- Some Pitfalls with OS Abstractions [if time]

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.43

High-Level vs. Low-Level File API

High-Level Operation:

```
size_t fread(...) {  
    Do some work like a normal fn...
```

```
asm code ... syscall # into %eax  
put args into registers %ebx, ...  
special trap instruction
```

```
Kernel:  
get args from regs  
dispatch to system func  
Do the work to read from the file  
Store return value in %eax
```

```
get return values from regs  
Do some more work like a normal fn...  
};
```

Low-Level Operation:

```
ssize_t read(...) {
```

```
asm code ... syscall # into %eax  
put args into registers %ebx, ...  
special trap instruction
```

```
Kernel:  
get args from regs  
dispatch to system func  
Do the work to read from the file  
Store return value in %eax
```

```
get return values from regs  
};
```

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.44

High-Level vs. Low-Level File API

- Streams are buffered in user memory:

```
printf("Beginning of line ");
sleep(10); // sleep for 10 seconds
printf("and end of line\n");
```

Prints out everything at once

- Operations on file descriptors are visible immediately

```
write(STDOUT_FILENO, "Beginning of line ", 18);
sleep(10);
write("and end of line \n", 16);
```

Outputs "Beginning of line" 10 seconds earlier than "and end of line"

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.45

What's in a FILE?

- What's in the FILE* returned by fopen?

- File descriptor (from call to open) <= Need this to interface with the kernel!
- Buffer (array)
- Lock (in case multiple threads use the FILE concurrently)

- Of course there's other stuff in a FILE too...

- ... but this is useful model to have

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.46

FILE Buffering

- When you call fwrite, what happens to the data you provided?
 - It gets written to the FILE's buffer
 - If the FILE's buffer is full, then it is *flushed*
 - » Which means it's written to the underlying file descriptor
 - The C standard library *may* flush the FILE more frequently
 - » e.g., if it sees a certain character in the stream
- When you write code, make the weakest possible assumptions about how data is flushed from FILE buffers

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.47

Example

```
char x = 'c';
FILE* f1 = fopen("file.txt", "w");
fwrite("b", sizeof(char), 1, f1);
FILE* f2 = fopen("file.txt", "r");
fread(&x, sizeof(char), 1, f2);
```

- The call to fread might see the latest write 'b'
- Or it might miss it and see end of file (in which case x will remain 'c')

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.48

Example

```
char x = 'c';
FILE* f1 = fopen("file.txt", "wb");
fwrite("b", sizeof(char), 1, f1);
fflush(f1);
FILE* f2 = fopen("file.txt", "rb");
fread(&x, sizeof(char), 1, f2);
```

- Now, the call to fread will definitely see the latest write 'b'

9/9/20

Kubiawicz CS162 © UCB Fall 2020

Lec 4.49

Writing Correct Code with FILE

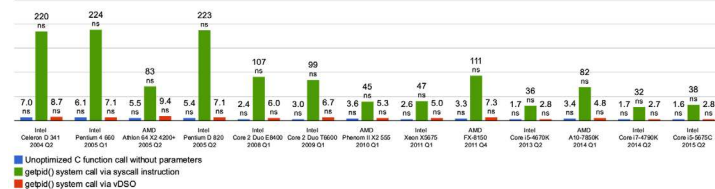
- Your code should behave correctly regardless of when C Standard Library flushes its buffer
 - Add your own calls to fflush so that data is written when you need to
 - Calls to fclose flush the buffer before deallocating memory and closing the file descriptor
- With the low-level file API, we don't have this problem
 - After write completes, data is visible to any subsequent reads

9/9/20

Kubiawicz CS162 © UCB Fall 2020

Lec 4.50

Why Buffer in Userspace? Overhead!



- Syscalls are 25x more expensive than function calls (~100 ns)
 - This example about special shared-memory interface to the getpid() functionality, but point is the same!
- **read/write** a file byte by byte? Max throughput of **~10MB/second**
- With **fgetc**? Keeps up with your SSD

9/9/20

Kubiawicz CS162 © UCB Fall 2020

Lec 4.51

Why Buffer in Userspace? Functionality!

- System call operations less capable
 - Simplifies operating system
- Example: No “read until new line” operation in kernel
 - Why? Kernel *agnostic* about formatting!
 - Solution: Make a big read syscall, find first new line in userspace
 - » i.e. use one of the following high-level options:

```
char *fgets(char *s, int size, FILE *stream);
ssize_t getline(char **lineptr, size_t *n, FILE *stream);
```

9/9/20

Kubiawicz CS162 © UCB Fall 2020

Lec 4.52

Today: The File Abstraction

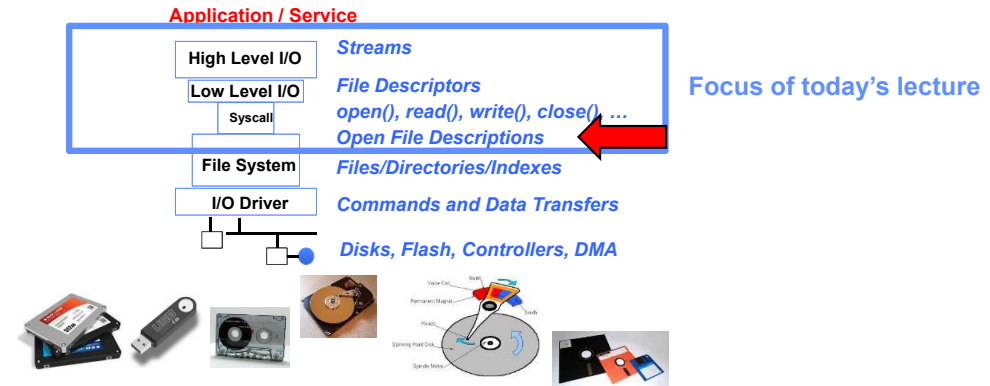
- High-Level File I/O: Streams
- Low-Level File I/O: File Descriptors
- *How* and *Why* of High-Level File I/O
- **Process State for File Descriptors**
- Some Pitfalls with OS Abstractions [if time]

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.53

I/O and Storage Layers



9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.54

State Maintained by the Kernel

- Recall: On a successful call to `open()`:
 - A file descriptor (int) is returned to the user
 - An open file description is created in the kernel
- For each process, kernel maintains mapping from file descriptor to open file description
 - On future system calls (e.g., `read()`), kernel looks up **open file description** using **file descriptor** and uses it to service the system call:

```
char buffer1[100];
char buffer2[100];
int fd = open("foo.txt", O_RDONLY);
read(fd, buffer1, 100);
read(fd, buffer2, 100);
```

The kernel remembers that the int it receives (stored in fd) corresponds to foo.txt

The kernel picks up where it left off in the file

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.55

What's in an Open File Description?

For our purposes, the two most important things are:

- Where to find the file data on disk
- The current position within the file

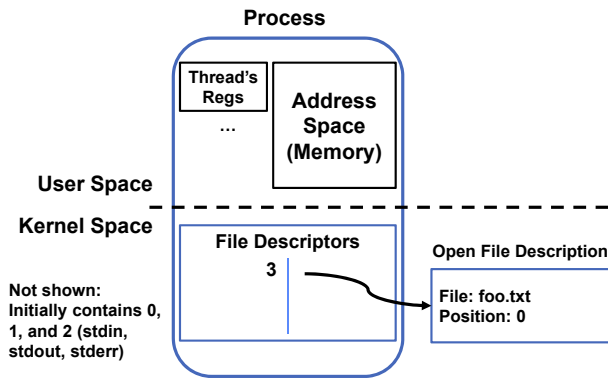
```
746 struct file {
747     union {
748         struct list_node fu_llist;
749         struct rcu_head fu_rcuhead;
750     } fu;
751     struct path f_path;
752     struct fdentry f_dentry;
753     struct inode f_inode;
754     const struct file_operations *f_op; /* caci */
755     /*
756      * Protects f_op.links, f_flags.
757      * Must not be taken from IRQ context.
758      */
759     spinlock_t f_lock;
760     atomic_long_t f_count;
761     unsigned int f_flags;
762     fmode_t f_mode;
763     struct mutex f_pos_lock;
764     loff_t f_pos;
765     struct fown_struct f_owner;
766     const struct cred *f_cred;
767     struct file_ra_state f_ra;
768     /*
769      * Used by fs/eventpoll.c to link all the hooks
770      *
771      * Used by fs/eventpoll.c to link all the hooks
772      */
773     void *f_security;
774     void *private_data;
775     /* needed for tty driver, and maybe others */
776     /*
777      * Used by fs/eventpoll.c to link all the hooks
778      */
779     struct list_head f_ep_links;
780     struct file_llink f_file_llink;
781     /* #ifdef CONFIG_EPOLL */
782     struct address_space *f_mapping;
783     /* #endif */
784     __attribute__((aligned(4))) /* test something weird
785     ...
```

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.56

Abstract Representation of a Process



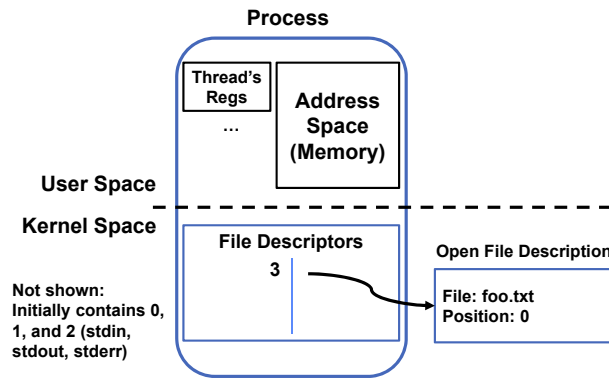
Suppose that we execute `open("foo.txt")` and that the result is 3

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.57

Abstract Representation of a Process



Suppose that we execute `open("foo.txt")` and that the result is 3

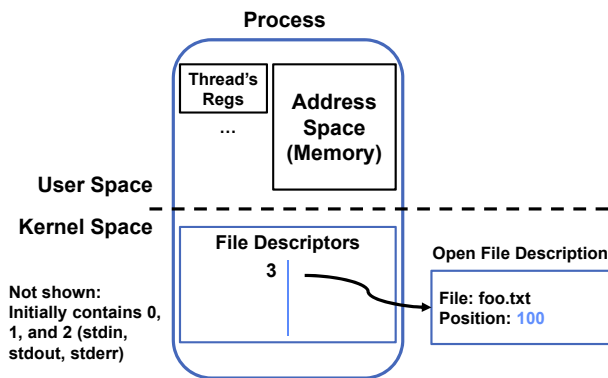
Next, suppose that we execute `read(3, buf, 100)` and that the result is 100

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.58

Abstract Representation of a Process



Suppose that we execute `open("foo.txt")` and that the result is 3

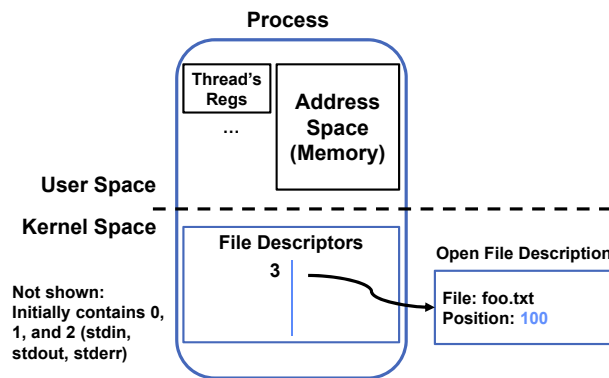
Next, suppose that we execute `read(3, buf, 100)` and that the result is 100

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.59

Abstract Representation of a Process



Suppose that we execute `open("foo.txt")` and that the result is 3

Next, suppose that we execute `read(3, buf, 100)` and that the result is 100

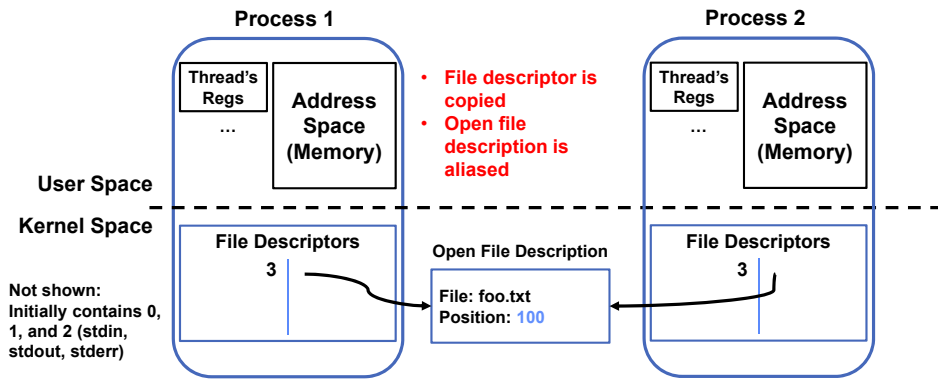
Finally, suppose that we execute `close(3)`

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.60

Instead of Closing, let's fork()!

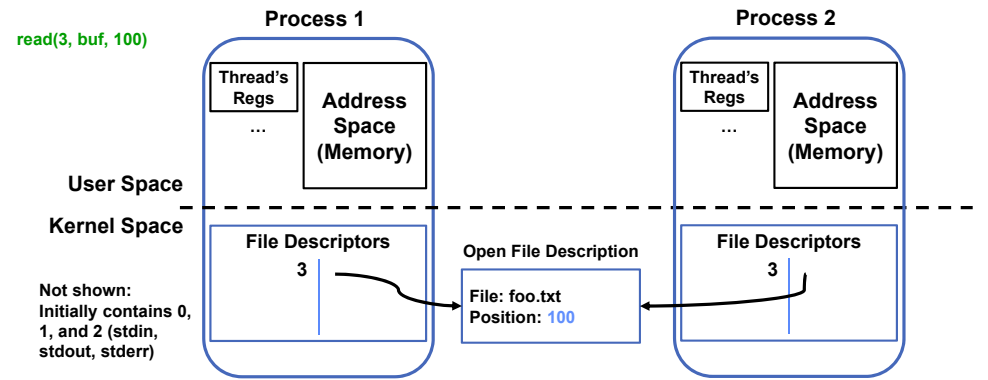


9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.61

Open File Description is Aliased

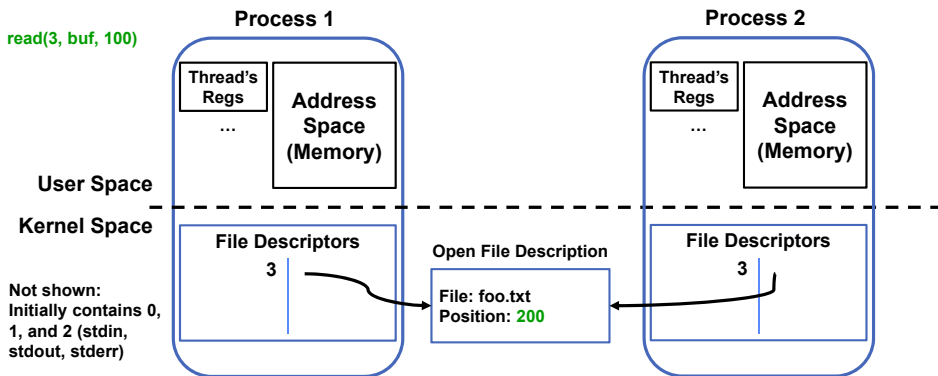


9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.62

Open File Description is Aliased

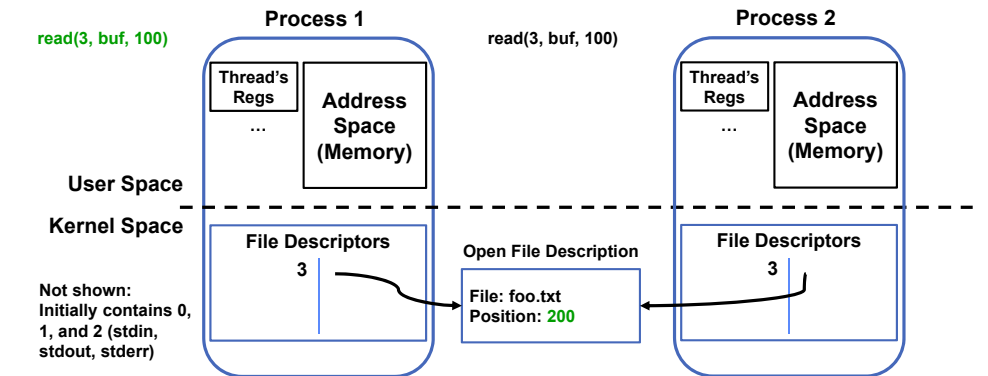


9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.63

Open File Description is Aliased

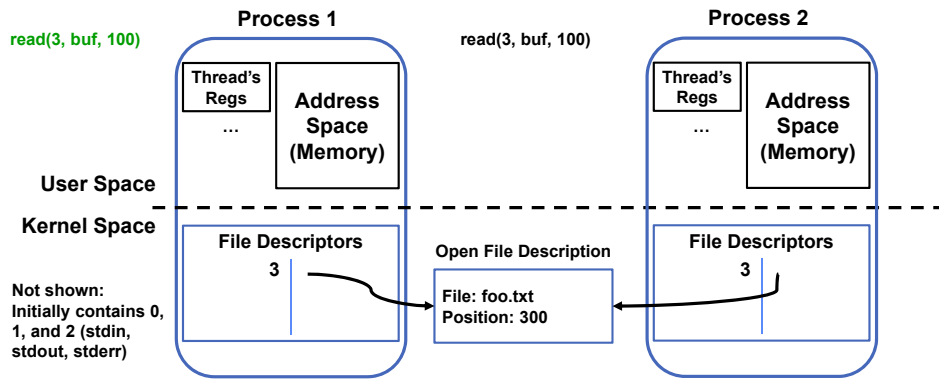


9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.64

Open File Description is Aliased

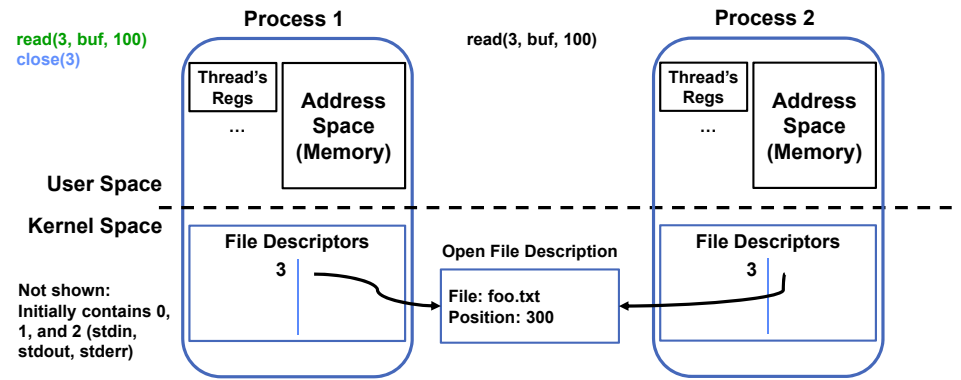


9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.65

File Descriptor is Copied

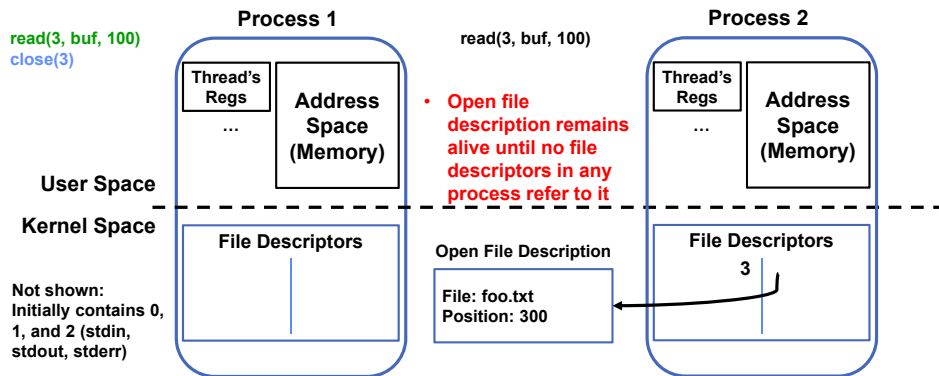


9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.66

File Descriptor is Copied



9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.67

Why is Aliasing the Open File Description a Good Idea?

- It allows for *shared resources* between processes

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.68

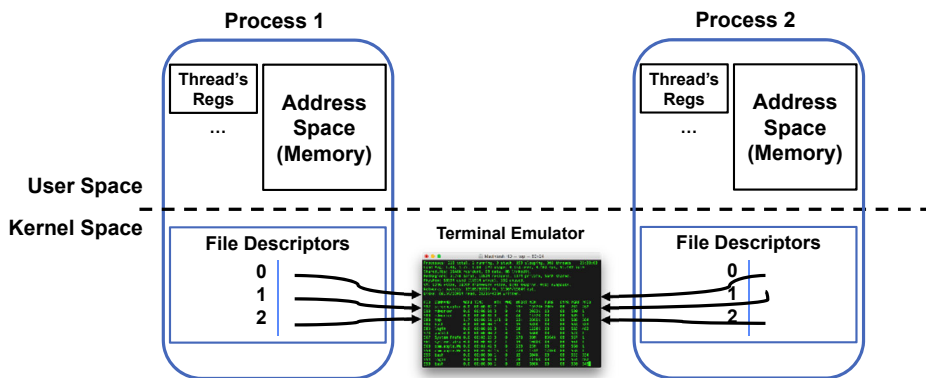
Recall: In POSIX, Everything is a “File”

- Identical interface for:
 - Files on disk
 - Devices (terminals, printers, etc.)
 - Regular files on disk
 - Networking (sockets)
 - Local interprocess communication (pipes, sockets)
- Based on the system calls **open()**, **read()**, **write()**, and **close()**

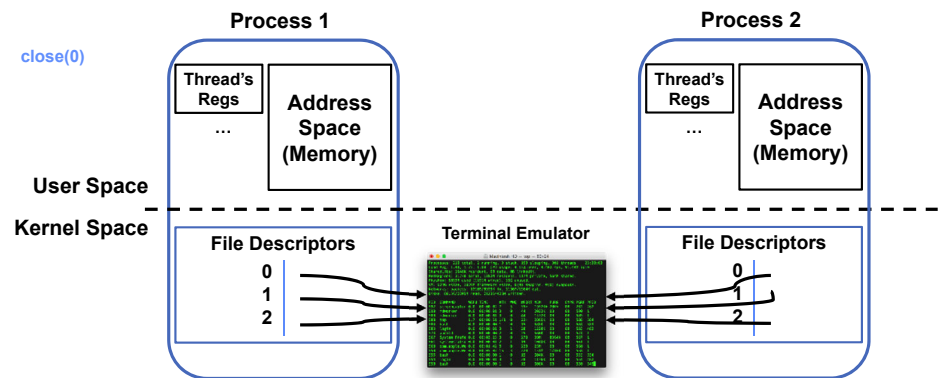
Example: Shared Terminal Emulator

- When you `fork()` a process, the parent’s and child’s `printf` outputs go to the same terminal

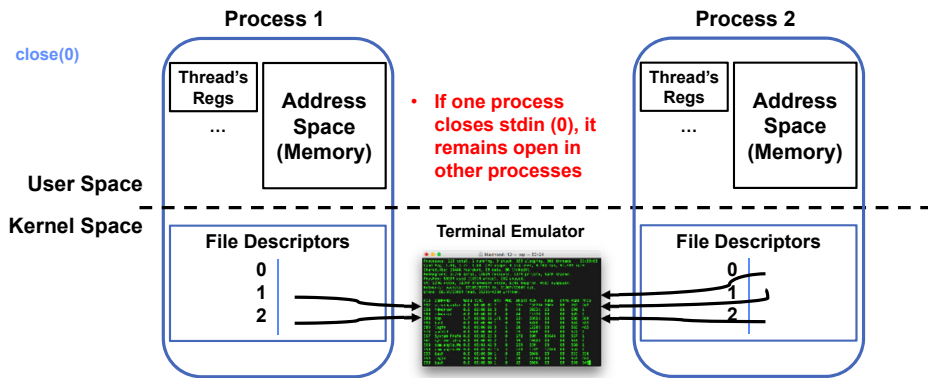
Example: Shared Terminal Emulator



Example: Shared Terminal Emulator



Example: Shared Terminal Emulator



9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.73

Other Examples

- Shared network connections after `fork()`
 - Allows handling each connection in a separate process
 - We'll explore this next time
- Shared access to pipes
 - Useful for interprocess communication
 - And in writing a shell (Homework 2)

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.74

Other Syscalls: dup and dup2

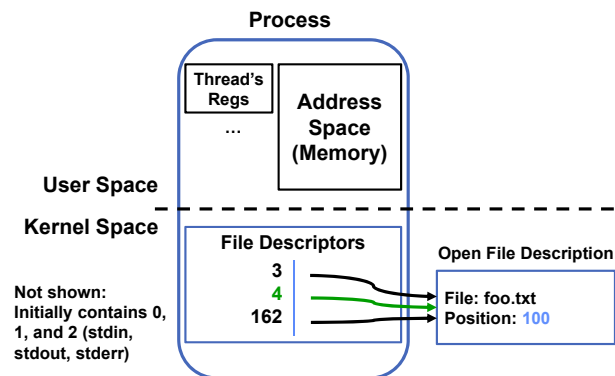
- They allow you to duplicate the file descriptor
- But the open file description remains aliased

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.75

Other Syscalls: dup and dup2



Suppose that we execute `open("foo.txt")` and that the result is 3

Next, suppose that we execute `read(3, buf, 100)` and that the result is 100

Next, suppose that we execute `dup(3)` and that the result is 4

Finally, suppose that we execute `dup2(3, 162)`

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.76

Today: The File Abstraction

- High-Level File I/O: Streams
- Low-Level File I/O: File Descriptors
- *How* and *Why* of High-Level File I/O
- Process State for File Descriptors
- **Some Pitfalls with OS Abstractions [if time]**

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.77

Unless you plan to call `exec()` in the child process

DON'T FORK() IN A PROCESS THAT ALREADY HAS MULTIPLE THREADS

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.78

fork() in Multithreaded Processes

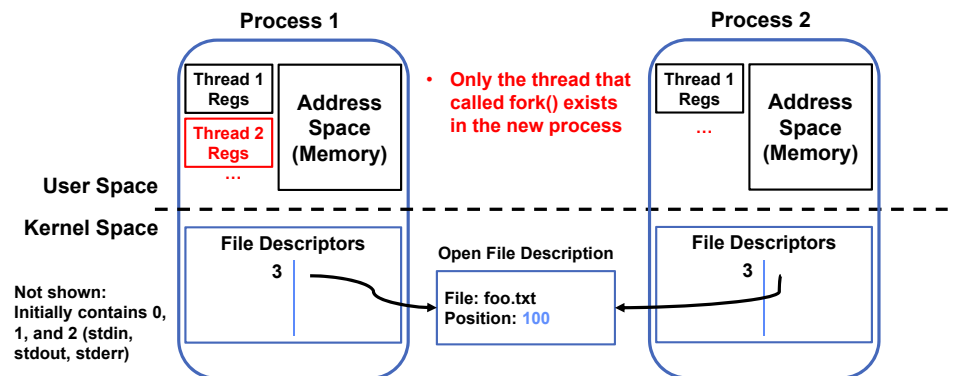
- The child process always has just a single thread
 - The thread in which `fork()` was called
- The other threads just vanish

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.79

fork() in a Multithreaded Processes



9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.80

Possible Problems with Multithreaded fork()

- When you call `fork()` in a multithreaded process, the other threads (the ones that didn't call `fork()`) just vanish
 - What if one of these threads was holding a lock?
 - What if one of these threads was in the middle of modifying a data structure?
 - No cleanup happens!
- It's safe if you call `exec()` in the child
 - Replacing the entire address space

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.81

**DON'T CARELESSLY MIX LOW-LEVEL
AND HIGH-LEVEL FILE I/O**

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.82

Avoid Mixing FILE* and File Descriptors

```
char x[10];
char y[10];
FILE* f = fopen("foo.txt", "rb");
int fd = fileno(f);
fread(x, 10, 1, f); // read 10 bytes from f
read(fd, y, 10); // assumes that this returns data starting at offset 10
```

- Which bytes from the file are read into `y`?
 - Bytes 0 to 9
 - Bytes 10 to 19
 - None of these?
- Answer: C! None of the above.
 - The `fread()` reads a big chunk of file into user-level buffer
 - Might be all of the file!

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.83

BE CAREFUL WITH FORK() AND FILE*

9/9/20

Kubiatowicz CS162 © UCB Fall 2020

Lec 4.84

Be Careful Using fork() with FILE*

```
FILE* f = fopen("foo.txt", "w");
fwrite("a", 1, 1, f);
fork();
fclose(f);
```

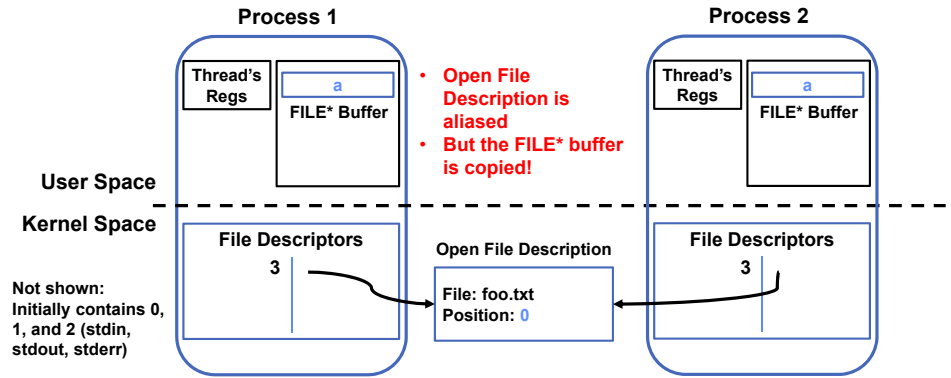
- Depends on whether this fwrite() call flushes...

After all processes exit, what is in foo.txt?

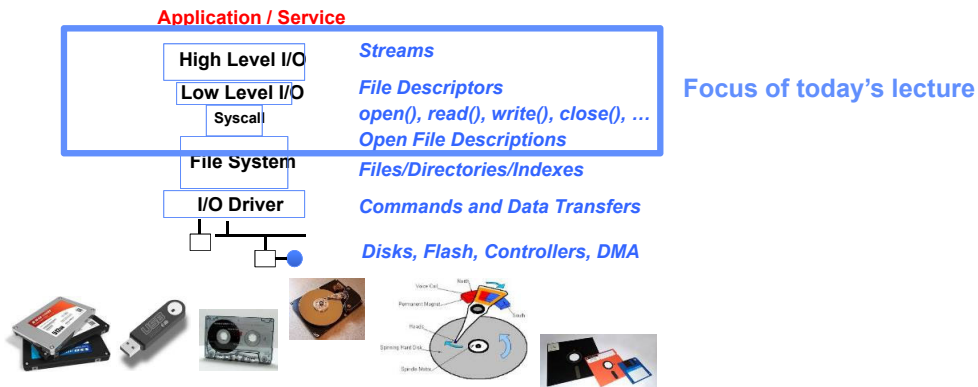
Could be either a or aa

- Usually aa based on what I've observed in Linux...

Be Careful Using fork() with FILE*



Conclusion



Conclusion

- POSIX idea: "everything is a file"
- All sorts of I/O managed by open/read/write/close
- We added two new elements to the PCB:
 - Mapping from file descriptor to open file description
 - Current working directory