

1.3 A Simple OS

- **OS Structure**
- Processes, Address Spaces, & Threads
- Managing Processes
- Loading Program Into Processes
- Files



Copyright © William C. Cheng



A Simple OS

- Sixth-Edition Unix
 - source license available to universities in 1975 from Bell Labs
 - had major influence on modern OSes
 - Solaris
 - Linux
 - MacOS X
 - Windows
- Fits into 64KB of memory
- single executable, completely stored in a **single file**
- **loaded** into memory as the OS boots
- **monolithic OS**



Copyright © William C. Cheng



User vs. Privileged Modes

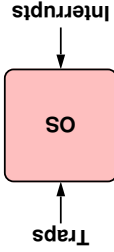
- **Processor modes:** part of the processor state (recall from your computer organization/architecture class regarding "processor")
 - most computers have at least two modes of execution
 - **user mode:** fewest privileges
 - **privileged mode:** most privileges
 - the only code that runs in this mode is part of the OS
- For Sixth-Edition Unix
 - the whole OS run in the privileged mode
 - everything else is an application and run in the user mode
- For other systems
 - major subsystems providing OS functionality may run in the user mode
- We use the word "**kernel**" to mean the portion of the OS that runs in privileged mode
- sometimes, a subset of this



Copyright © William C. Cheng



A Simple OS Structure

- Application programs call upon the OS via **traps**
 - External devices call upon the OS via **interrupts**
- 

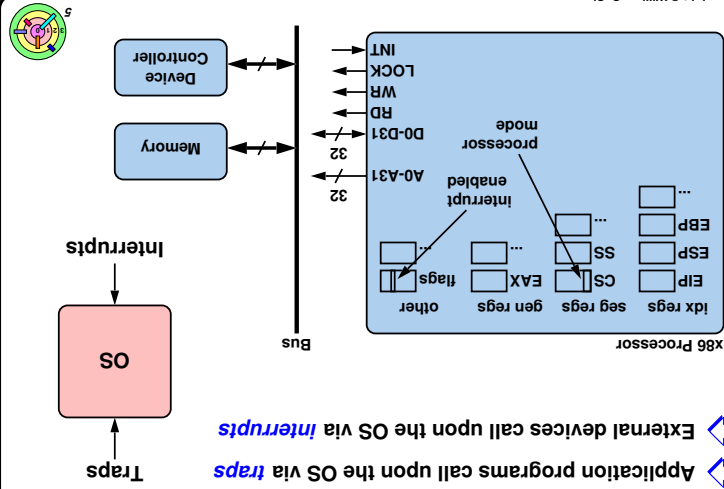


Copyright © William C. Cheng



A Simple OS Structure

- Application programs call upon the OS via **traps**
- External devices call upon the OS via **interrupts**



Copyright © William C. Cheng



Traps

- **Traps** are the general means for invoking the kernel from user code
- although we usually think of traps as **errors**
 - divide by zero, segmentation fault, bus error, etc.
 - but they don't have to be
 - **system calls**, **page fault**, etc.
- Traps always elicit some sort of response
 - for programming errors, the default action is to **terminate**
 - the user program
 - for system calls, the OS is asked to perform some service



Copyright © William C. Cheng





Copyright © William C. Cheng

- ### A Special Kind Of Trap - System Calls
- Invoking OS functionality in the kernel is more complex
 - but we want to make it look simple to applications
 - must be done carefully and correctly
 - really cannot trust the application programmers to do the right thing every time
 - Provide **system calls** through which user code can access the kernel **in a controlled manner**
 - any necessary checking on whether the request should be permitted can be done in the system call
 - all done in user mode
 - it all goes well
 - sets things up
 - traps** into the kernel by executing a special machine instruction, i.e., the "trap" machine instruction
 - the kernel figures out why it was invoked and handles the trap

Operating Systems - CSCI 402



Copyright © William C. Cheng

- ### Upcall
- A program may establish a handler (i.e., a **signal handler**) to be invoked in response to the error
 - the handler might clean up after the error and then terminate the program, or it might perform corrective action and continue with normal execution
 - The **upcall** mechanism
 - signals** allow the kernel to invoke code that's part of user program
 - for example, you can set a timer to expire at a certain time, when it expires, the OS can use the upcall mechanism to call a specified user function

Operating Systems - CSCI 402



Copyright © William C. Cheng

- ### Program Execution
- Fundamental **abstraction of program execution**
 - memory
 - address space
 - things that are addressable by the program are kept together here
 - in Sixth-Edition Unix, processes do not share address space
 - recall that **process** is an abstraction of **memory**
 - processor(s)
 - recall that **thread** is an abstraction of **processor**
 - "execution context"
 - the **state** of a process and its threads
 - exactly "where you are" in the program
 - a thread needs some sort of a context to execute
 - Note: multiple meanings of the word "context" in this class
 - save context** and **restore context**
 - thread context** vs. **interrupt context**

Operating Systems - CSCI 402



Copyright © William C. Cheng

- ### Interrupts
- An **interrupt** is a request from an **external device** for a response from the **processor**
 - handled independently of any user program
 - unlike a trap, which is handled as part of the program that caused the trap
 - response to a trap directly affects that program
 - response to an interrupt may or may not indirectly affect the currently running program
 - often has **no direct effect** on the currently running program
 - There's also something called **software interrupt**
 - generated programmatically (i.e., not by a device) by executing an "interrupt" machine instruction
 - this is very different from a hardware interrupt, although the mechanisms of handling interrupts are all very similar

Operating Systems - CSCI 402



Copyright © William C. Cheng

- OS Structure
- Processes, Address Spaces, & Threads**
- Managing Processes
- Loading Program Into Processes
- Files

1.3 A Simple OS

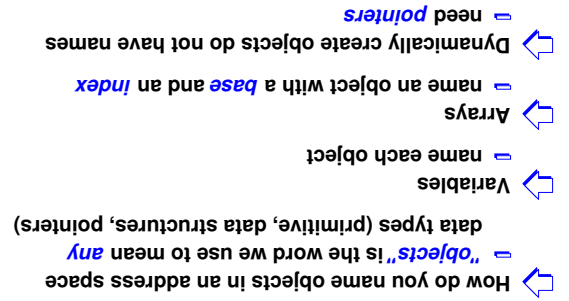
Operating Systems - CSCI 402



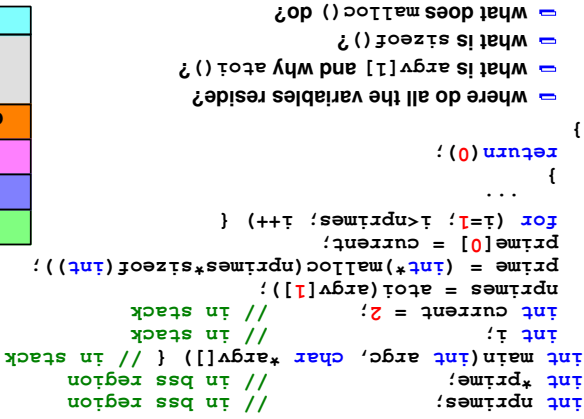
Copyright © William C. Cheng

- ### Program Execution
- With abstraction, comes an interface / API
 - for processes
 - `fork()`, `exec()`, `wait()`, `exit()`

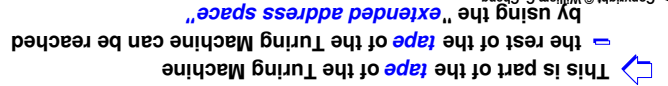
Operating Systems - CSCI 402



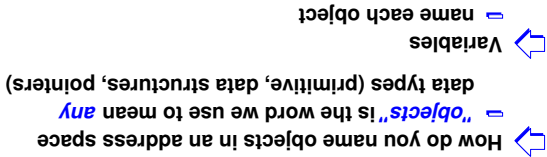
Operating Systems - CSCI 402



Operating Systems - CSCI 402



Operating Systems - CSCI 402 -



Operating Systems - CSCI 402

My color codes for code

Operating Systems - CSCI 402 -

-

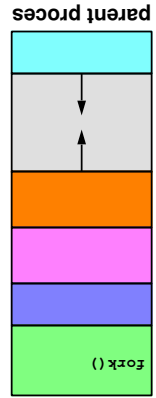
Operating Systems - CSCI 402

1.3 A Simple OS

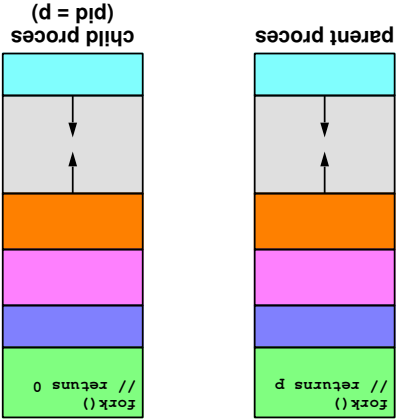
- OS Structure
- Processes, Address Spaces, & Threads
- Managing Processes
- Loading Program Into Processes
- Files



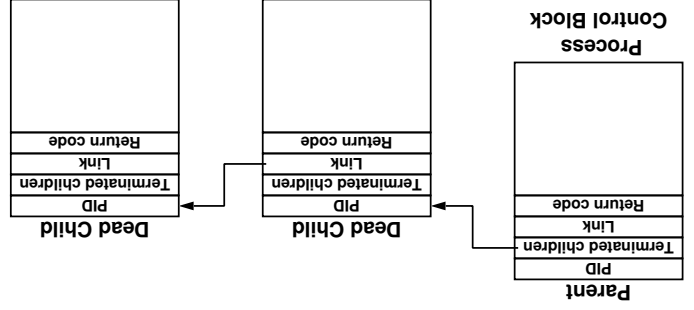
- Creating a process is deceptively simple
 - make a copy of a process (the parent process)
 - pid_t fork(void)
 - the process where fork() is called is the **parent** process
 - the copy is the **child** process
 - in a way, fork() returns twice
 - once in the parent, the returned value is the **process ID (PID)** of the child process
 - once in the child, the returned value is 0
 - a PID is 16-bit long
 - this is the **only** way to create a process
- Making a copy of the entire address space can be expensive
 - Ch 7 shows speed up tricks
 - e.g., text segment is read-only so parent and child can share it
- Example: relationship between a shell (i.e., a command interpreter, such as /bin/tcsh) and /bin/ls



Creating a Process: Before



Creating a Process: After



Process Control Blocks

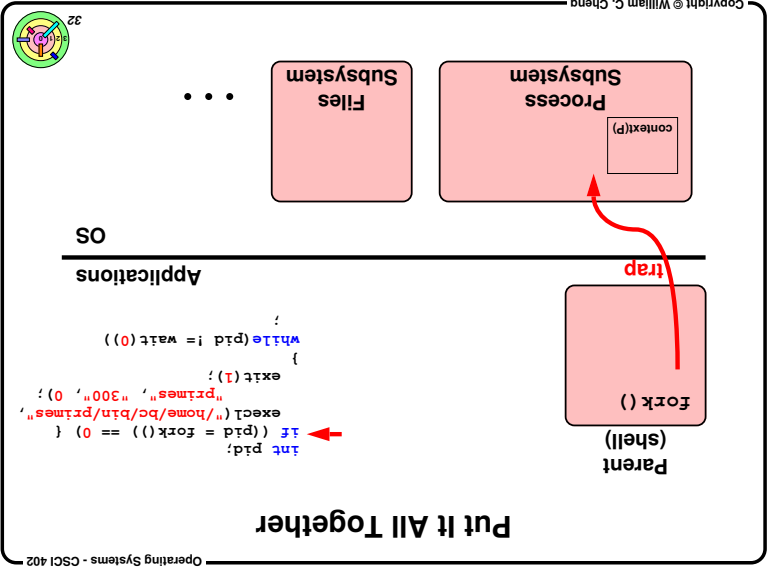
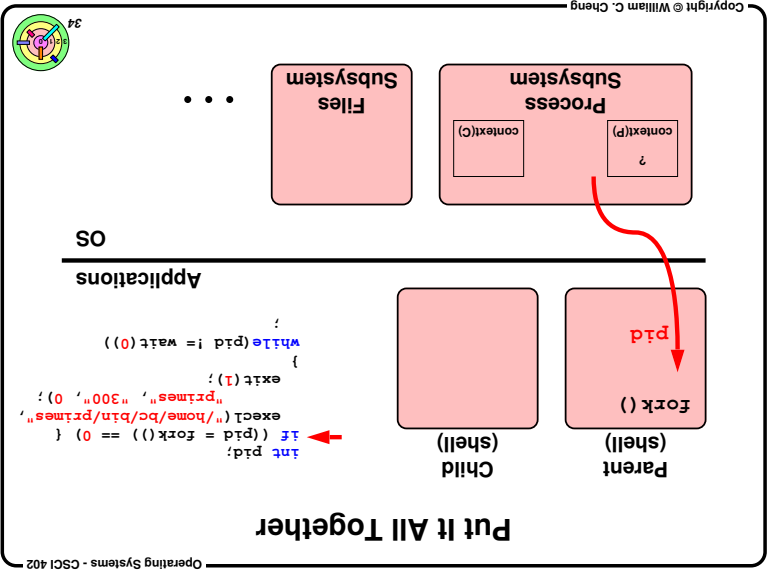
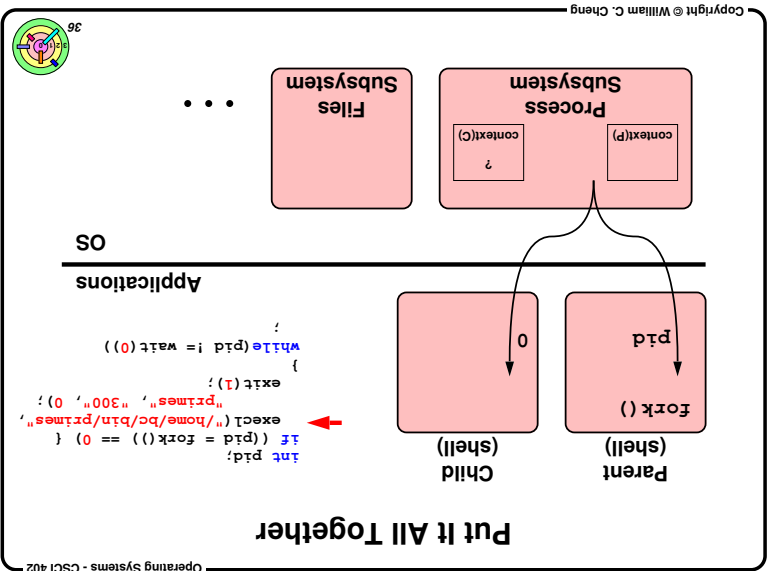
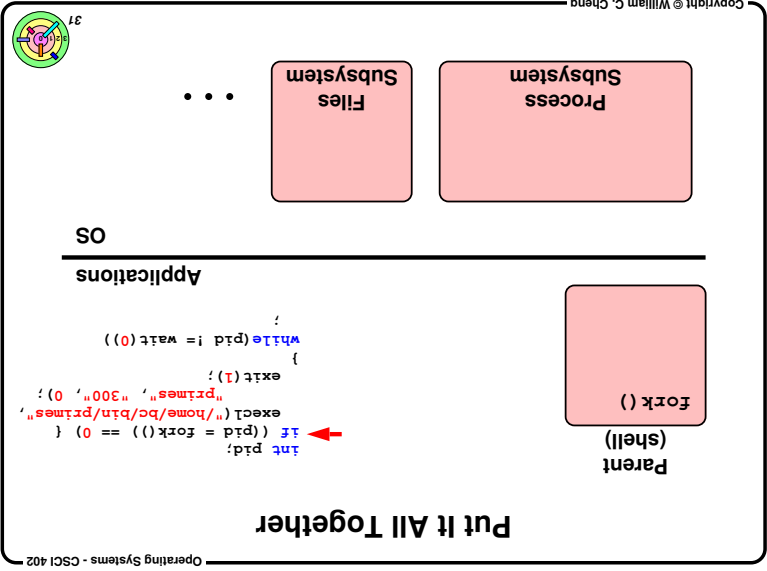
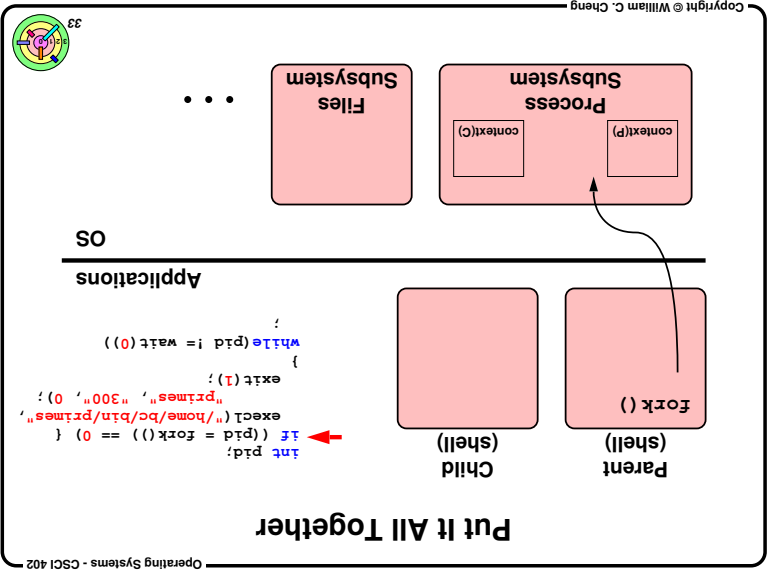
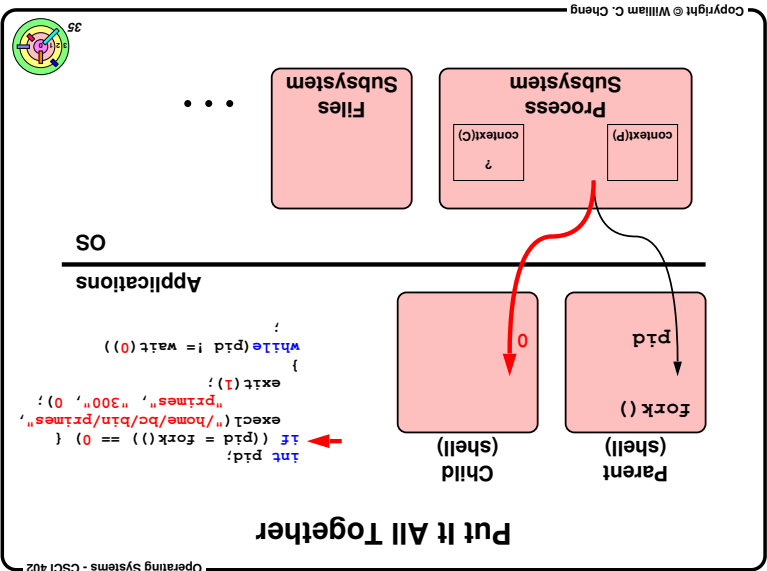


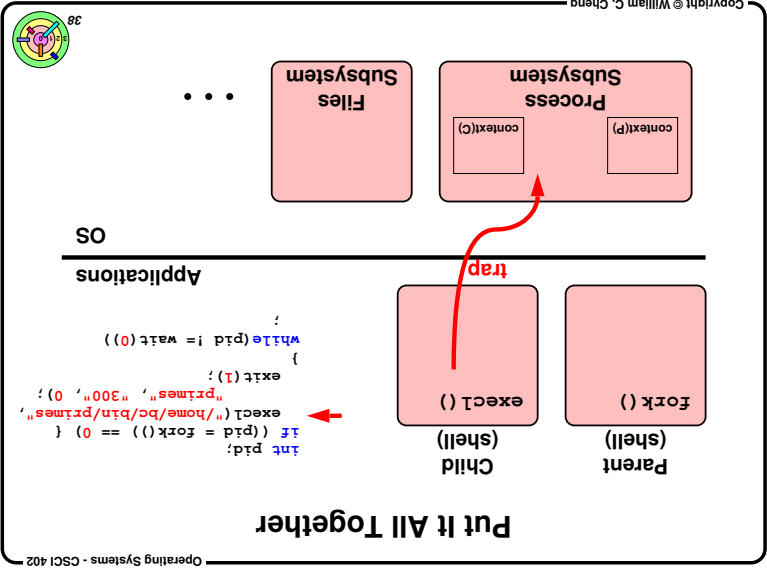
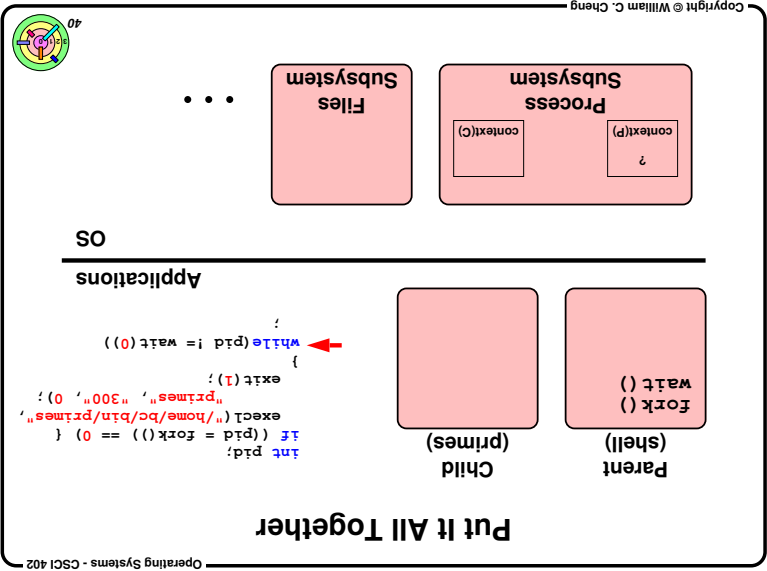
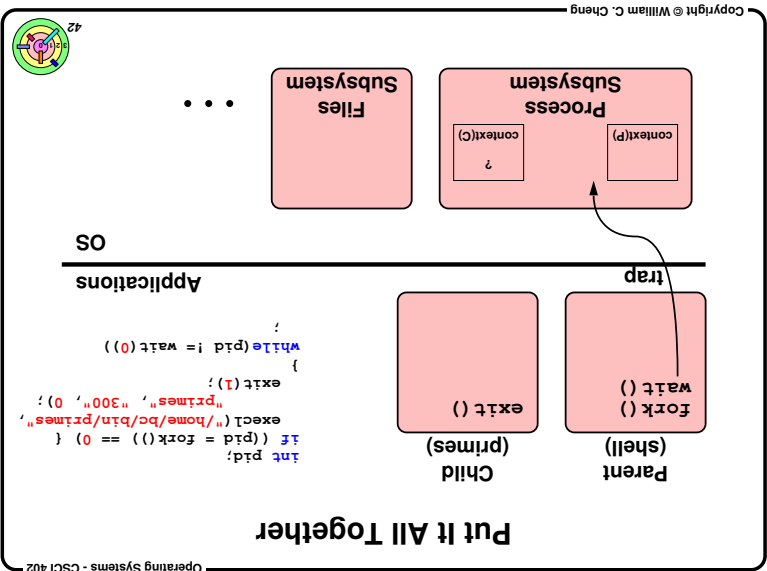
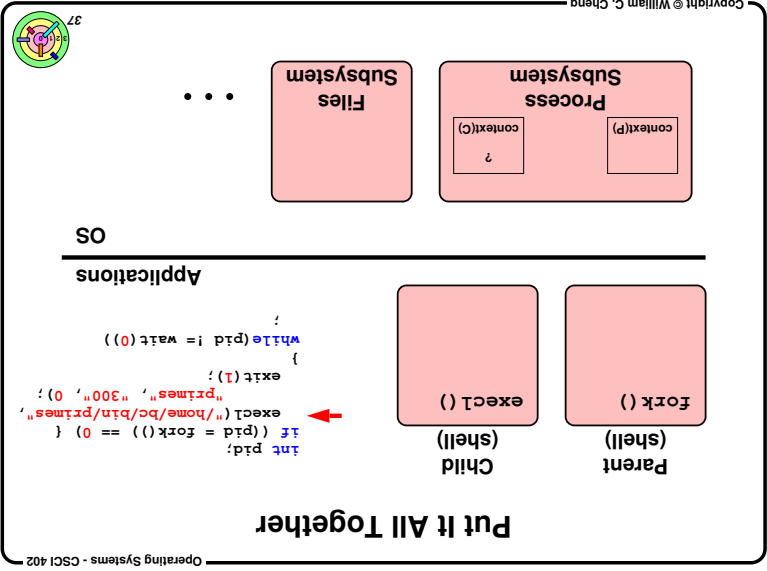
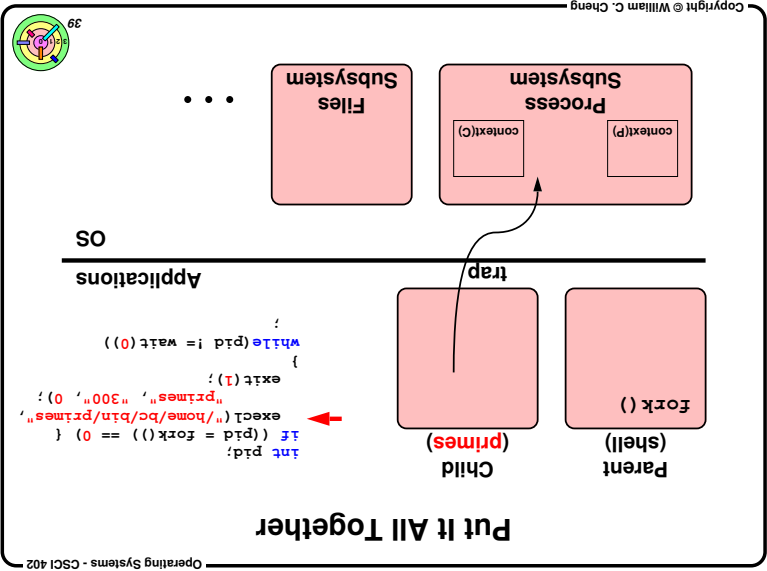
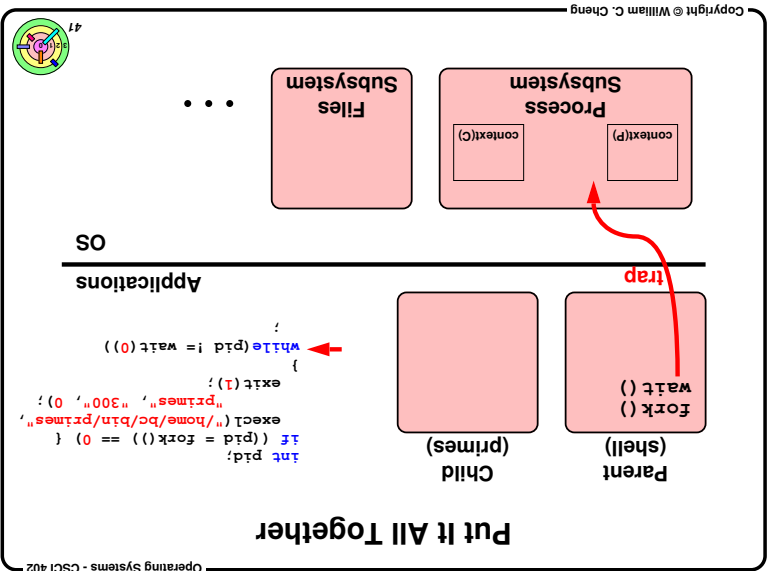
- Process Control Block (PCB) is a kernel data structure
- pretty much every field is unsigned
- return code (when a process dies) is 8-bit long
- so that the parent process can know what happened to child
- the "Link" field points to the next PCB
- but, the next PCB in what list?

Fork and Wait

```
short pid;
if ((pid = fork()) == 0) {
    /* some code is here for the child to execute */
    exit(n);
} else {
    int ReturnCode;
    while (pid != wait(&ReturnCode))
        /* the child has terminated with ReturnCode as
        its return code */
        ;
    e.g., /bin/tcsh forks /bin/ls
    what does exit(n) do other than copying n into PCB?
    least significant 8-bits of n
    what happens when main() calls return(n)?
    eventually, exit(n) will be invoked
    pid_t wait(int *status) is a blocking call
    it reaps dead child processes one at a time
}
```

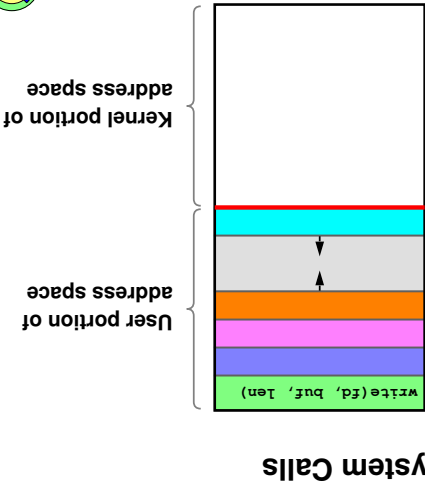




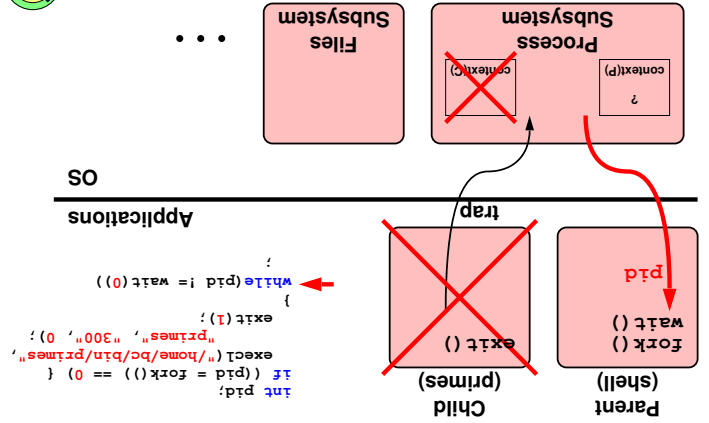


Operating Systems - CSCI 402

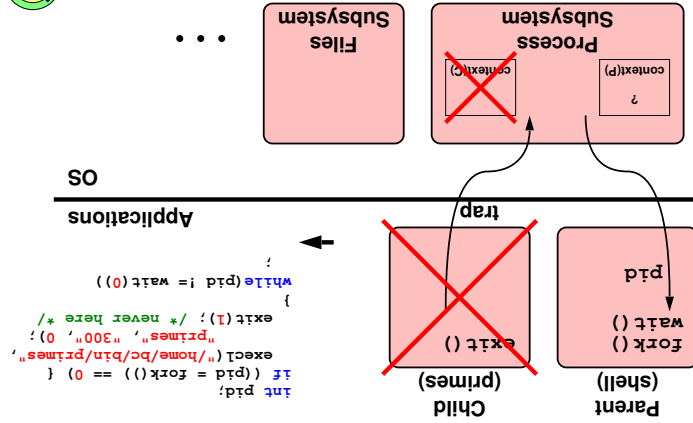
Operating Systems - CSCI 402



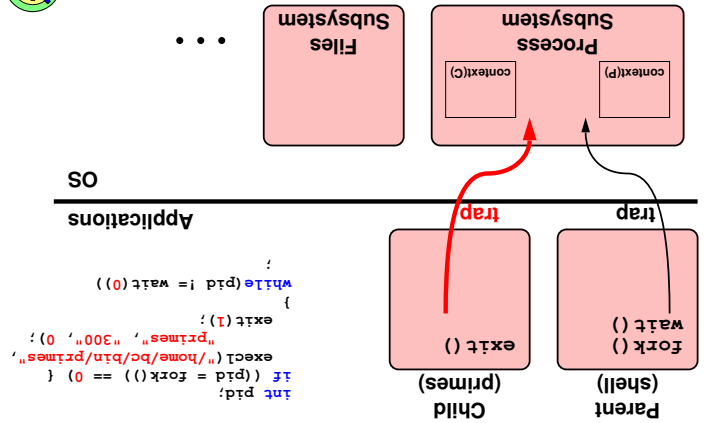
Operating Systems - CSCI 402



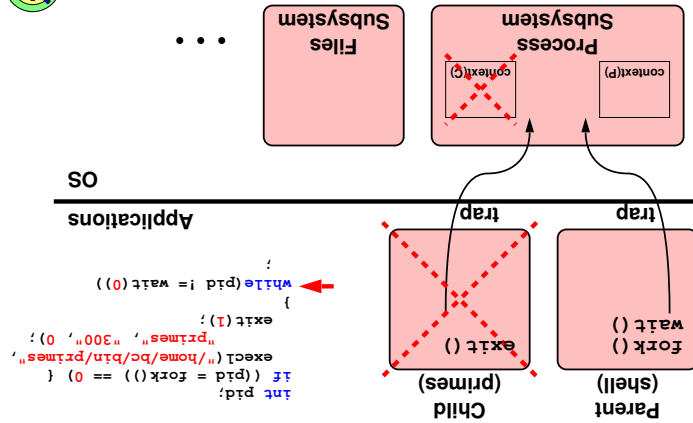
Operating Systems - CSCI 402



Operating Systems - CSCI 402



Operating Systems - CSCI 402



- ## Naming Files

Operating Systems - CSCI 402

- ## 1.3 A Simple OS

Operating Systems - CSCI 402

- ## System Calls

Operating Systems - CSCI 402

Copyright © William C. Cheng

Allocation of File Descriptors

Whenever a process requests a new file descriptor, the lowest numbered file descriptor not already associated with an open file is selected; thus

```
#include <fcntl.h>
#include <unistd.h>
...
close(0);
fd = open("file", O_RDONLY);
```

will always associate "file" with file descriptor 0 (assuming that open() succeeds)

You will need to implement the above rule in the kernel 2 assignment

Operating Systems - CSCI 402

Copyright © William C. Cheng

Back to Primes

Have our primes program write out the solution, i.e., the primes [array

```
int nprimes;
int *primes;
int main(int argc, char *argv[]) {
    ...
    for (i=1; i<nprimes; i++) {
        if (write(1, primes, nprimes*sizeof(int)) == -1) {
            perror("primes output");
            exit(1);
        }
        return(0);
    }
```

the output is not readable by human

Operating Systems - CSCI 402

Copyright © William C. Cheng

File Handles (File Descriptors)

```
int fd;
char buffer[1024];
int count;
if ((fd = open("/home/bc/file", O_RDWR)) == -1) {
    // the file couldn't be opened
    perror("/home/bc/file");
    exit(1);
}
if ((count = read(fd, buffer, 1024)) == -1) {
    // the read failed
    perror("read");
    exit(1);
}
// buffer now contains count bytes read from the file
```

- what is O_RDWR?
- what does perror() do?
- cursor position in an opened file depends on what functions/system calls you use
- what about C++?

Operating Systems - CSCI 402

Copyright © William C. Cheng

Running It

```
if (fork() == 0) {
    // set up file descriptor 1 in the child process *
    close(1);
    if (open("/home/bc/output", O_WRONLY) == -1) {
        perror("/home/bc/output");
        exit(1);
    }
    execl("/home/bc/bin/primes", "primes", "300", "0");
    exit(1);
}
while (pid != wait(0)) /* ignore the return code */
/* parent continues here */
}
space
file descriptors are allocated lowest first on open()
extended address space survives execs
new code is same as running
% primes 300 > /home/bc/output
```

Operating Systems - CSCI 402

Copyright © William C. Cheng

Human-Readable Output

please see the *Programming FAQ* regarding the difference between a *file descriptor* and a *file pointer*

```
int nprimes;
int *primes;
int main(int argc, char *argv[]) {
    ...
    for (i=1; i<nprimes; i++) {
        for (j=0; j<nprimes; j++) {
            printf("%d\n", primes[j]);
        }
        return(0);
    }
```

Operating Systems - CSCI 402

Copyright © William C. Cheng

Standard File Descriptors

- 0 is stdin (by default, "map/connect" to the keyboard)
- 1 is stdout (by default, "map/connect" to the display)
- 2 is stderr (by default, "map/connect" to the display)

```
main() {
    char buf[BUFFSIZE];
    const char *note = "write failed\n";
    while ((n = read(0, buf, sizeof(buf))) > 0) {
        if (write(1, buf, n) != n) {
            void write(2, note, strlen(note));
            exit(EXIT_FAILURE);
        }
        return(EXIT_SUCCESS);
    }
```

Operating Systems - CSCI 402

-
- The diagram illustrates the kernel data structures for file sharing. It is divided into two main sections: **User** and **Kernel**.
- User Space:**
- address space**: Contains a **User** and a **File descriptor**.
- Kernel Space:**
- File-descriptor table (per process)**: A table with indices 0, 1, 2, 3, ..., n-1. An arrow points from the **File descriptor** in the user space to an entry in this table.
 - system file table (system-wide)**: A table with columns: **ref count**, **access mode**, **file location**, and **inode pointer**. An arrow points from an entry in the **File-descriptor table** to this table.
 - cursor**: Points to the **file location** column.
 - this is yet another pointer**: Points to the **inode pointer** column.

Operating Systems - CSCI 402

Operating Systems - CSCI 402

The ">" parameter in a shell command that instructs the command

- If ">" weren't there, the output would go to the display

Can also redirect input


```
% cat < /home/dc/output
- when the "cat" program reads from file descriptor 0, it would
```

```
get the data bytes from the file "/home/bc/Output"
```

Operating Systems - CSCI 402

- It also refers to the *process's* current *context* for that file
- includes how the file is to be accessed (how open () was

invoked
cursor position

 **Context** information must be maintained by the OS and not

directly by the user program

- later on it calls `write()` using the opened file descriptor
- how does the OS knows that it doesn't have write access?

- stores `0_RDONLY` in context

change O_RDONLY to O_RDWR

- all it can see is the handle

- the tile handle is an *index* into an array maintained for the process in kernel's address space

the process in kernel's address space