

Exceptional Control Flow: Exceptions and Processes

CSE4009: System Programming

Woong Sul

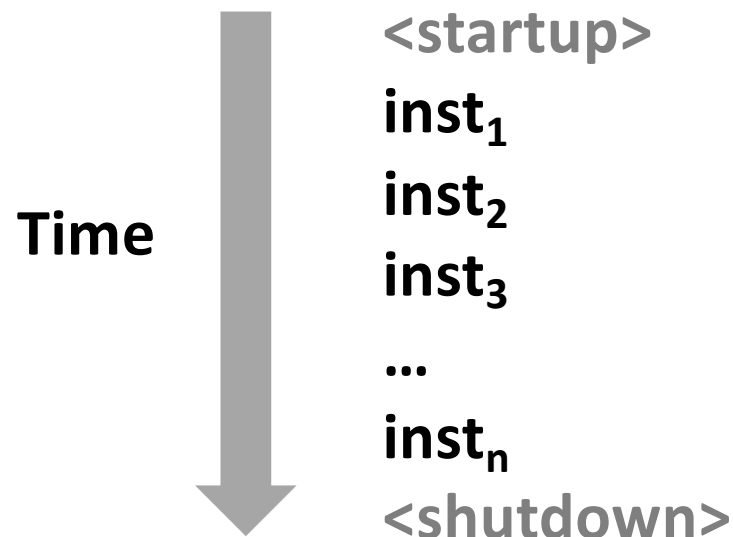
Today

- Exceptional Control Flow
- Exceptions
- Processes
- Process Control

Control Flow

- Processors do only one thing:
 - From startup to shutdown, a CPU simply reads and executes (interprets) a sequence of instructions, one at a time
 - This sequence is the CPU's control flow (or flow of control)

(Physical) control flow



Altering the Control Flow

- Up to now: two mechanisms for changing control flow:
 - Jumps and branches
 - Call and return→ React to changes in *program state*
- Insufficient for a useful system:
Difficult to react to changes in *system state*
 - Data arrives from a disk or a network adapter
 - Instruction divides by zero
 - User hits Ctrl-C at the keyboard
 - System timer expires
- System needs mechanisms for “**E**xceptional **C**ontrol **F**low”

Exceptional Control Flow

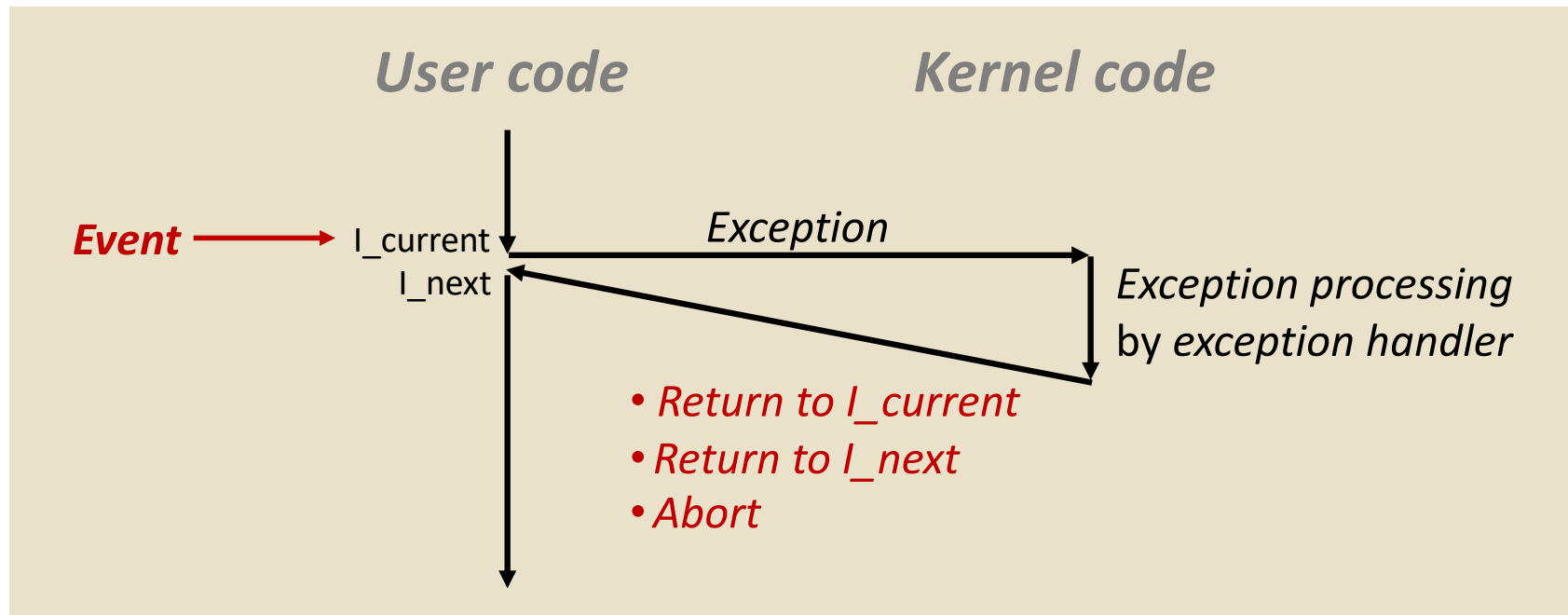
- Exists at all levels of a computer system
- Low level mechanisms
 1. **Exceptions**
 - Change in control flow in response to a system event (i.e., change in system state)
 - Implemented using combination of hardware and OS software
- Higher level mechanisms
 2. **Process context switch**
 - Implemented by OS software and hardware timer
 3. **Signals**
 - Implemented by OS software

Today

- Exceptional Control Flow
- **Exceptions**
- Processes
- Process Control

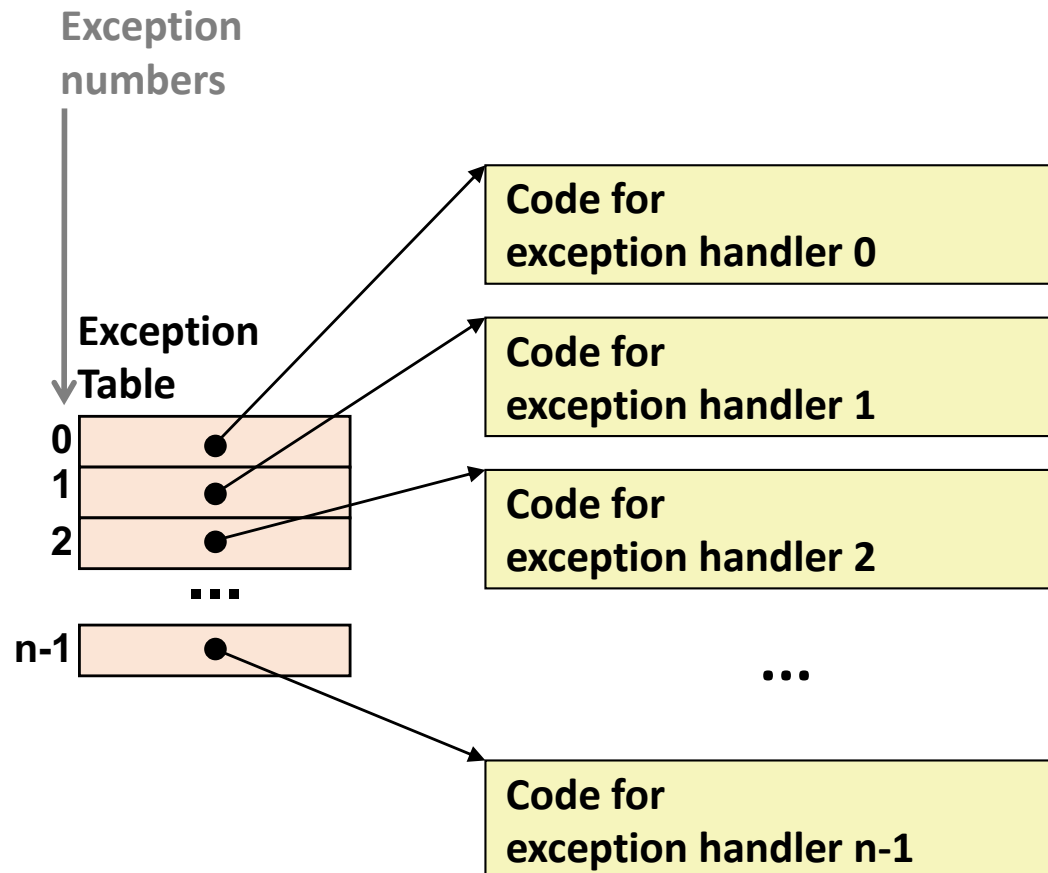
Exceptions

- An **exception** is a transfer of control to the OS kernel in response to some *event* (i.e., change in processor state)
 - Kernel is the memory-resident part of the OS
 - Examples of events: Divide by 0, arithmetic overflow, page fault, I/O request completes, typing Ctrl-C



Exception Tables

- Exception handlers might be implemented by SW



- Each type of event has a unique exception number k
- k = index into exception table (a.k.a. interrupt vector)
- Handler k is called each time exception k occurs

- How to call the handler?
 - i.e., How does packet arrival interrupt the current flow

Asynchronous Exceptions (Interrupts)

- Caused by events external to the processor
 - Indicated by setting the processor's *interrupt pin*
 - Handler returns to “next” instruction
- Examples:
 - Timer interrupt
 - Every few ms, an external timer chip triggers an interrupt
 - Used by the kernel to take back control from user programs
 - I/O interrupt from external device
 - Hitting Ctrl-C at the keyboard
 - Arrival of a packet from a network
 - Arrival of data from a disk

Synchronous Exceptions

- Caused by events that occur as a result of executing an instruction:
 - **Traps**
 - Intentional → like procedure calls
 - Examples: **system calls**, breakpoint traps, special instructions
 - Returns control to “next” instruction
 - **Faults**
 - Unintentional but possibly recoverable
 - Examples: page faults (recoverable), protection faults (unrecoverable), floating point exceptions
 - Either re-executes faulting (“current”) instruction or aborts
 - **Aborts**
 - Unintentional and unrecoverable
 - Examples: illegal instruction, parity error, machine check
 - Aborts current program

System Calls

- Each x86-64 system call has a unique ID number
- Examples:

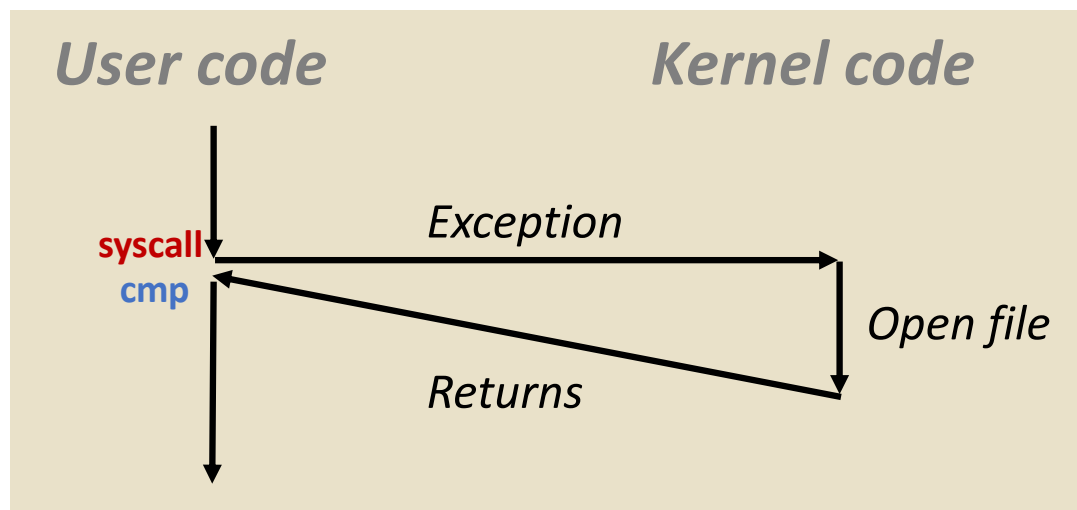
<i>Number</i>	<i>Name</i>	<i>Description</i>
0	read	Read file
1	write	Write file
2	open	Open file
3	close	Close file
4	stat	Get info about file
57	fork	Create process
59	execve	Execute a program
60	_exit	Terminate process
62	kill	Send signal to process

System Call Example: Opening File

- User calls: `open(filename, options)`
- Calls `__open()`, which invokes system call `syscall`

```

00000000000e5d70 <__open>:
...
e5d79:  b8 02 00 00 00      mov  $0x2,%eax  # open is syscall #2
e5d7e:  0f 05              syscall         # Return value in %rax
e5d80:  48 3d 01 f0 ff ff   cmp  $0xffffffffffffffff001,%rax
...
e5dfa:  c3                retq
    
```



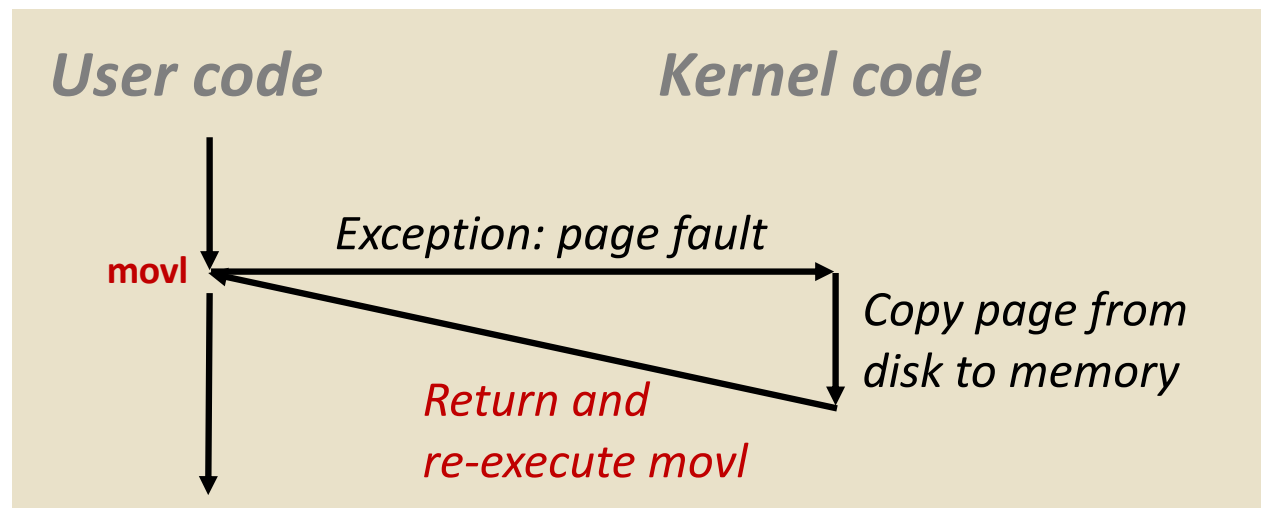
- `%rax` contains syscall number
- Other arguments in `%rdi`, `%rsi`, `%rdx`, `%r10`, `%r8`, `%r9`
- Return value in `%rax`
- Negative value is an error corresponding to negative `errno`

Fault Example: Page Fault

- User writes to memory location
- That portion (page) of user's memory is currently on disk

```
int a[1000];
main ()
{
    a[500] = 13;
}
```

80483b7: c7 05 10 9d 04 08 0d movl \$0xd,0x8049d10

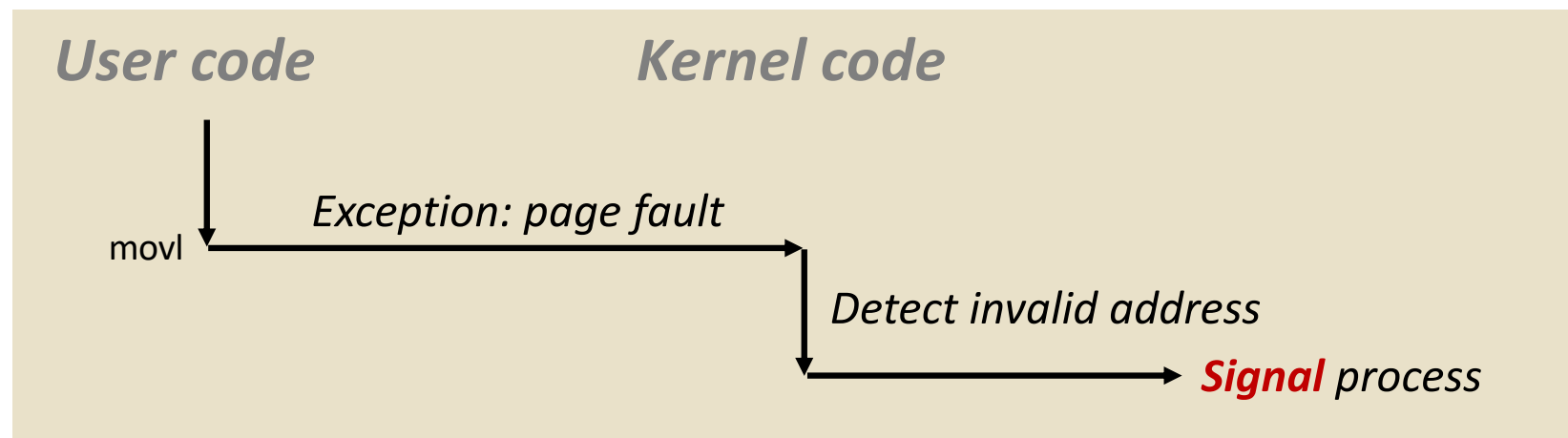


Fault Example: Invalid Memory Reference

- Sends SIGSEGV signal to user process
- User process exits with “segmentation fault”

```
int a[1000];
main ()
{
    a[5000] = 13;
}
```

80483b7:	c7 05 60 e3 04 08 0d	movl	\$0xd,0x804e360
----------	----------------------	------	-----------------

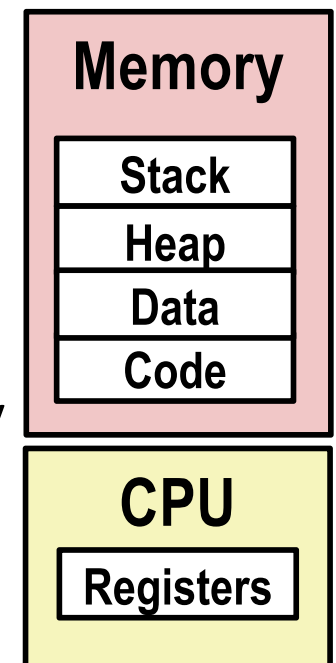


Today

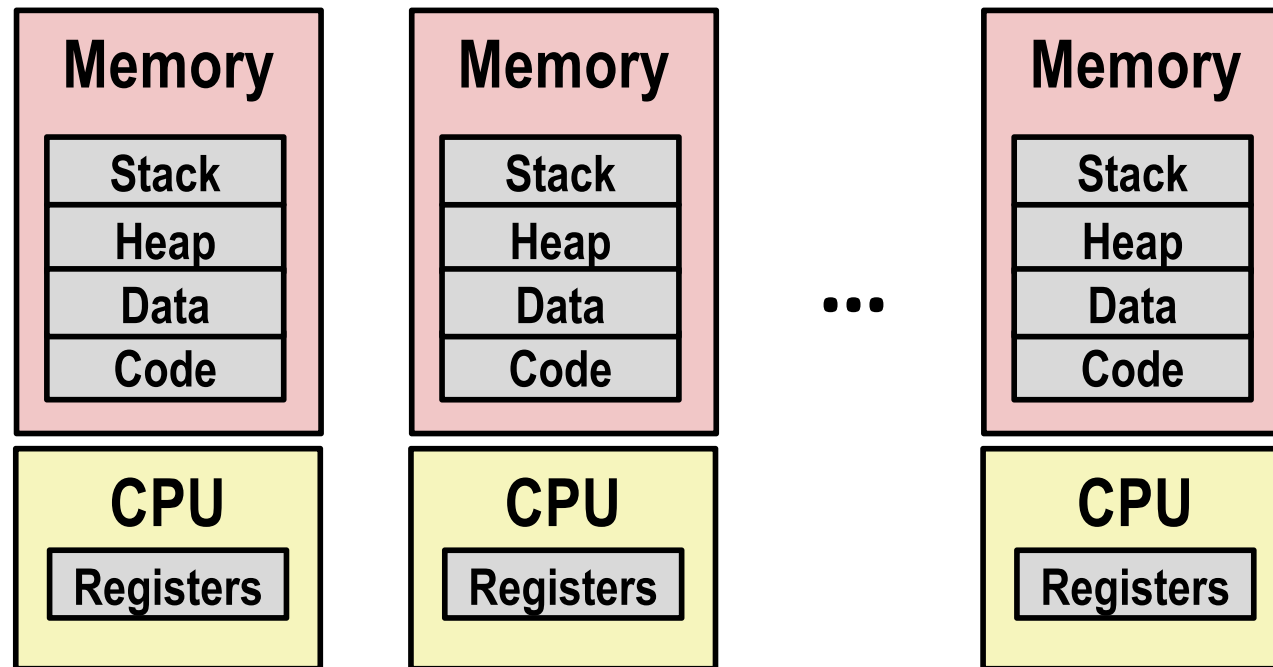
- Exceptional Control Flow
- Exceptions
- **Processes**
- Process Control

Processes

- An instance of a ***running program***
 - One of the most profound ideas in computer science
 - Not the same as “program” or “processor”
- Process provides each program with two key abstractions:
 - ***Logical control flow***
 - Each program seems to have exclusive use of the CPU
 - Provided by kernel mechanism called ***context switching***
 - ***Private address space***
 - Each program seems to have exclusive use of main memory
 - Provided by kernel mechanism called ***virtual memory***



Multiprocessing: The Illusion



- Computer runs many processes simultaneously
 - Applications for one or more users
 - Web browsers, email clients, editors, ...
 - Background tasks
 - Monitoring network & I/O devices

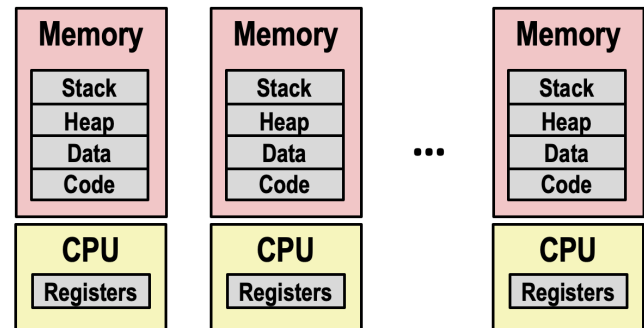
Multiprocessing Example

- Running program “top” on Mac
 - System has 737 processes, 2 of which are active
 - Identified by Process ID (PID)

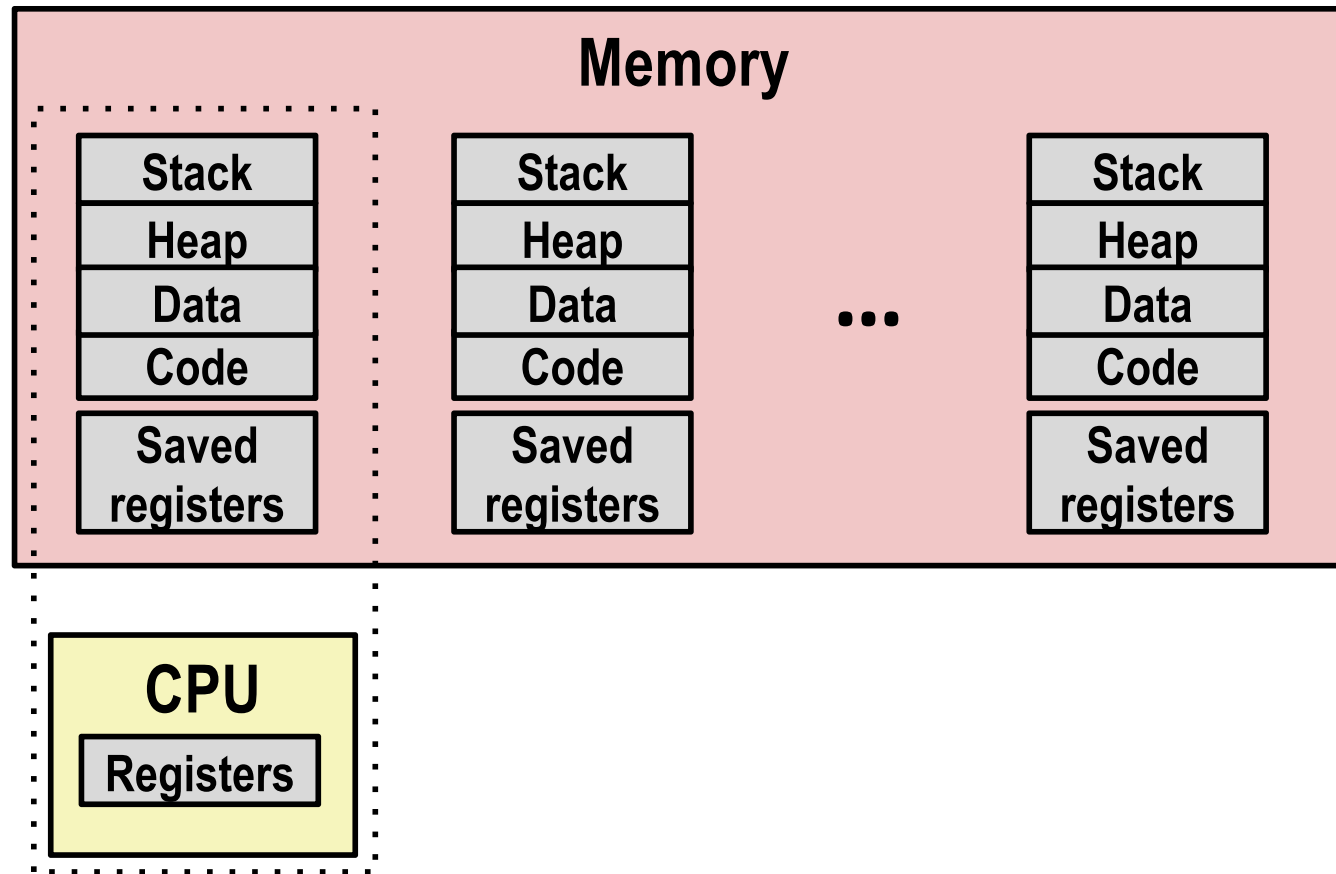
```

woongsul — woongsul@Woongs-MacBook-Pro: ~ — top — 80x24
~ — woongsul@Woongs-MacBook-Pro: ~ — top
Processes: 737 total, 2 running, 735 sleeping, 3017 threads      16:45:51
Load Avg: 1.06, 1.60, 1.84  CPU usage: 1.44% user, 1.67% sys, 96.88% idle
SharedLibs: 695M resident, 79M data, 42M linkedit.
MemRegions: 131891 total, 3752M resident, 244M private, 4433M shared.
PhysMem: 16G used (3372M wired, 771M compressor), 118M unused.
VM: 24T vsize, 4374M framework vsize, 9059629(64) swapins, 10017426(0) swapouts.
Networks: packets: 5939601/6938M in, 3048111/520M out.
Disks: 10032748/183G read, 3535442/92G written.

PID    COMMAND      %CPU  TIME    #TH    #WQ    #PORT  MEM    PURG    CMPRS  PGRP
27992   top           8.2   00:04.48 1/1     0      31     8440K  0B     0B     27992
172     WindowServer  7.4   03:00:53 22      13     5755+  1452M- 81M    157M- 172
27995   screencaptur  3.1   00:00.31 7        5     210+   14M+   0B     0B     27995
0       kernel_task   2.9   77:28.85 339/12  0      0      403M+  0B     0B     0
165     bluetoothd    1.8   21:32.93 11       7     362    10M    384K  29     165
439     aciseagend    1.7   03:16.21 6        1     62     8236K  0B     41     439
27994   screencaptur  1.6   00:00.20 12      11     76+    4592K+ 620K  0B     27994
270     airportd     1.2   51:58.77 18      16     319-   20M+   0B     10     270
441     TouchBarServ  1.0   15:49.26 7        3     393    26M    5888K  83     441
18827   Microsoft Po  0.9   78:11.37 51       6     5023-  1240M- 2415M  13     18827
26439   com.apple.Ap  0.8   00:52.78 4        3     392    2996K  0B     19     26439
544     nearbyd      0.7   07:18.69 8        6     99     2556K  0B     82     544
23635   Dropbox      0.5   02:13.42 148      1     821    332M   0B     103M   23635
3141   Terminal     0.5   15:15.23 11       4     584    86M    9364K  39M    3141
    
```

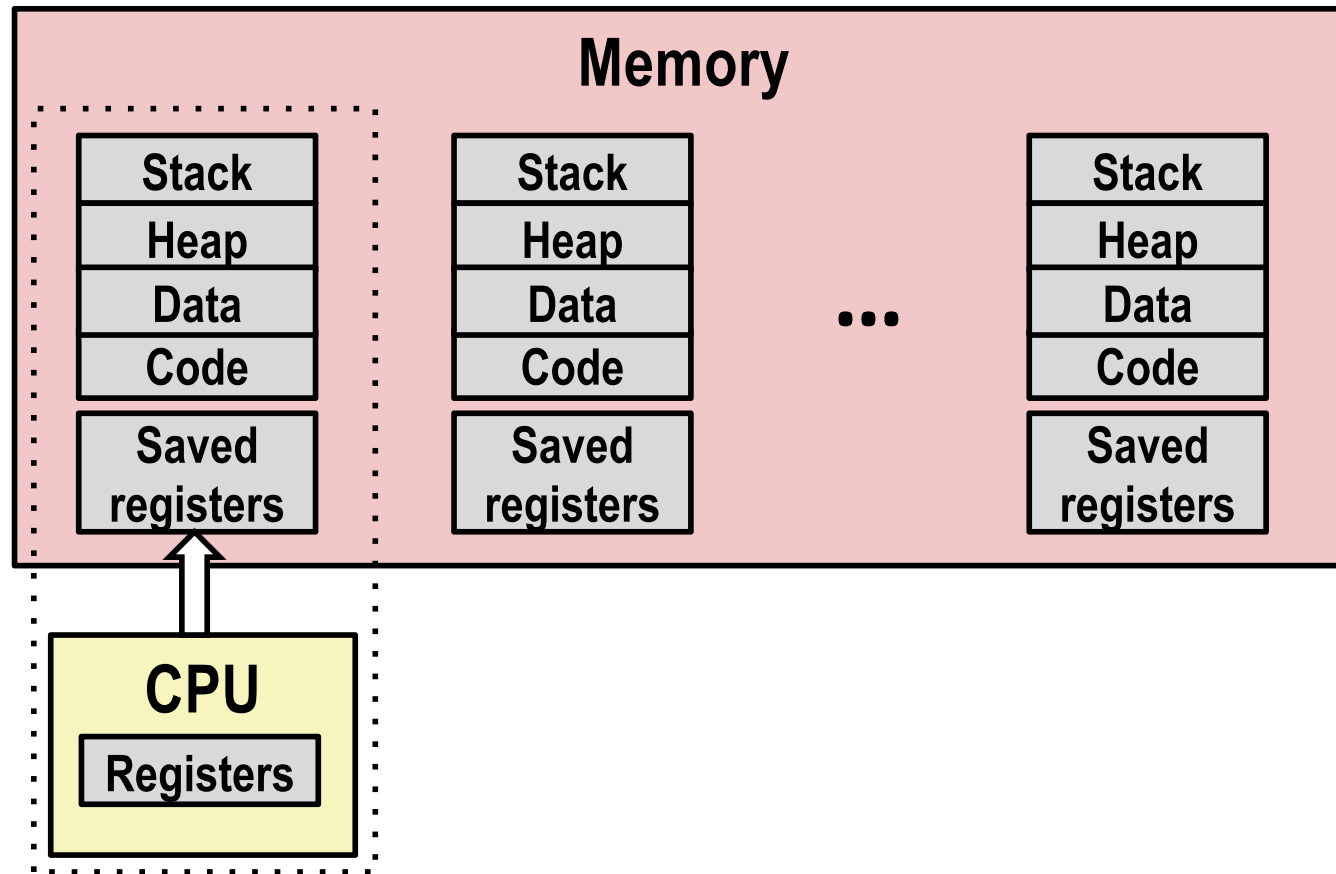


Multiprocessing: The (Traditional) Reality



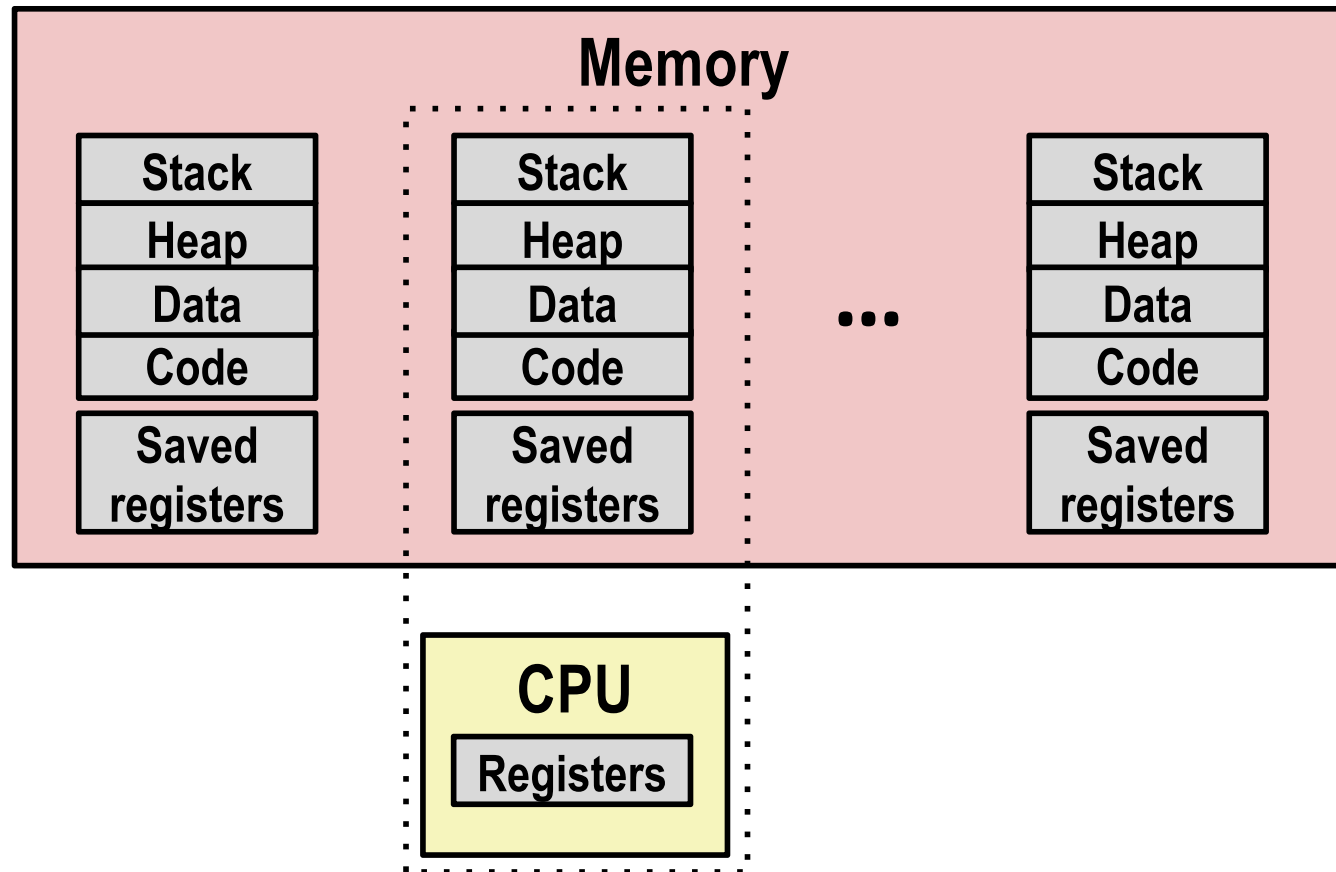
- Single processor executes multiple processes concurrently
 - Process executions interleaved (multitasking)
 - Address spaces managed by virtual memory system (later in course)
 - Register values for non-executing processes saved in memory

Multiprocessing: The (Traditional) Reality



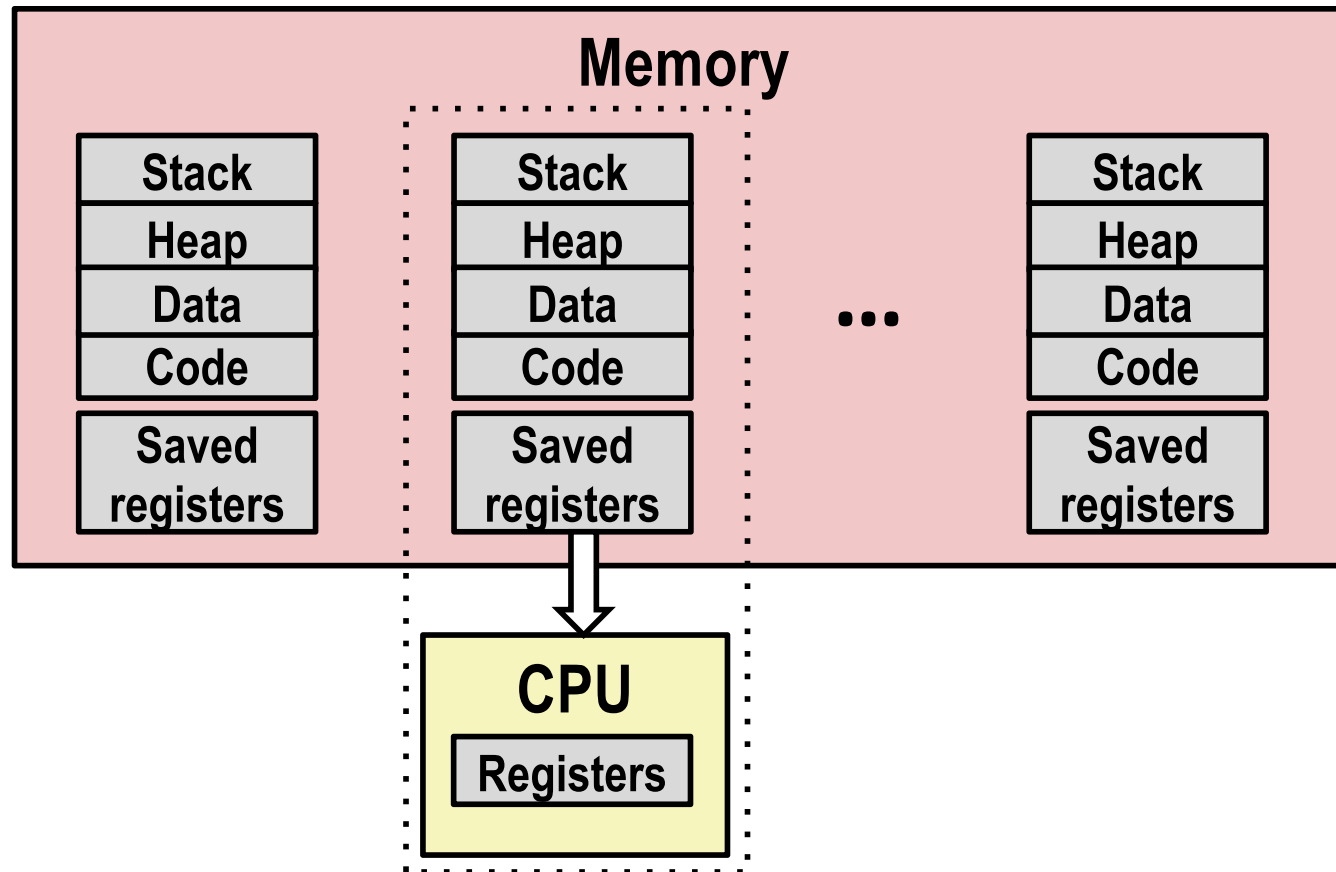
- Save current registers in memory

Multiprocessing: The (Traditional) Reality



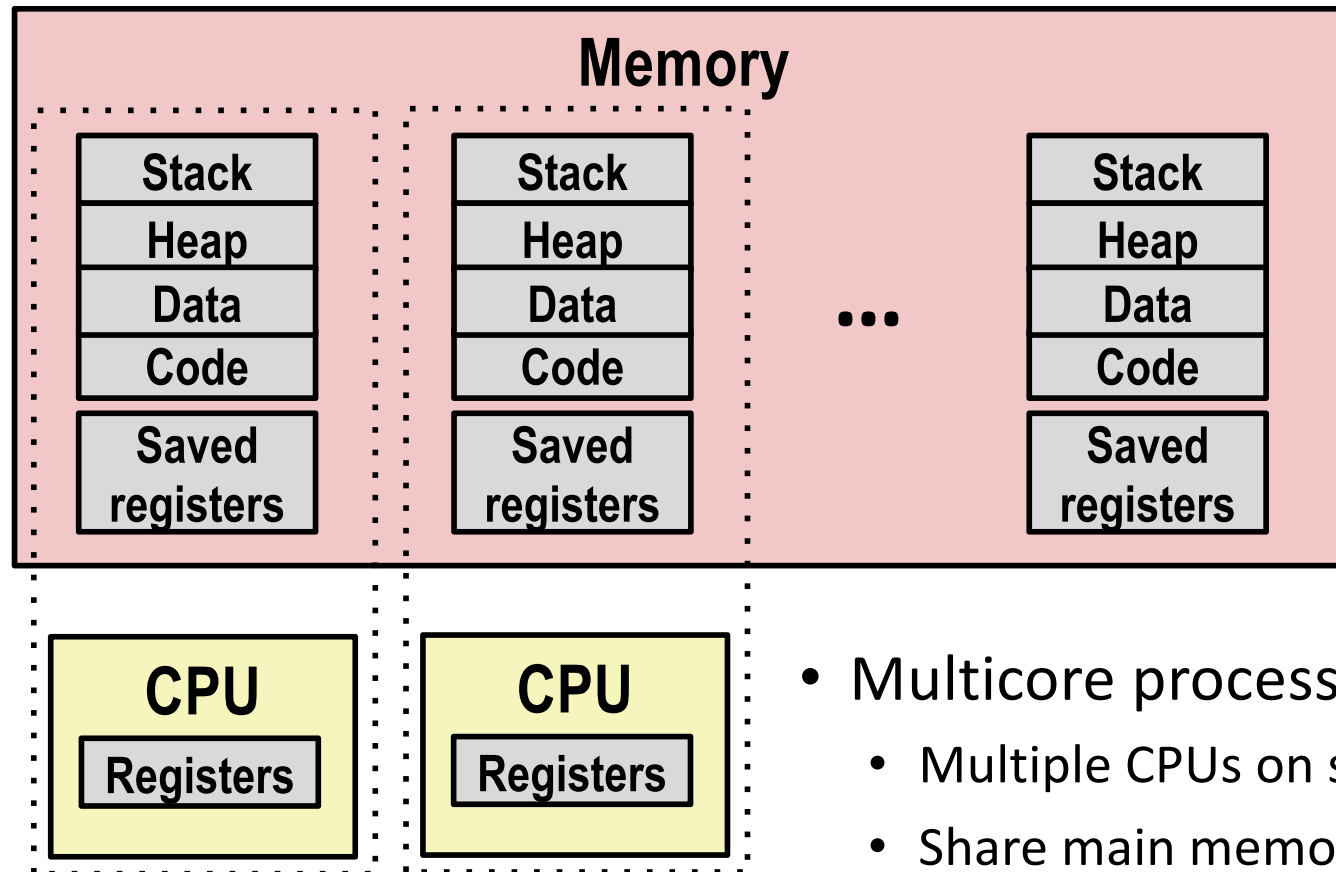
- Schedule next process for execution

Multiprocessing: The (Traditional) Reality



- Load saved registers and switch address space (context switch)

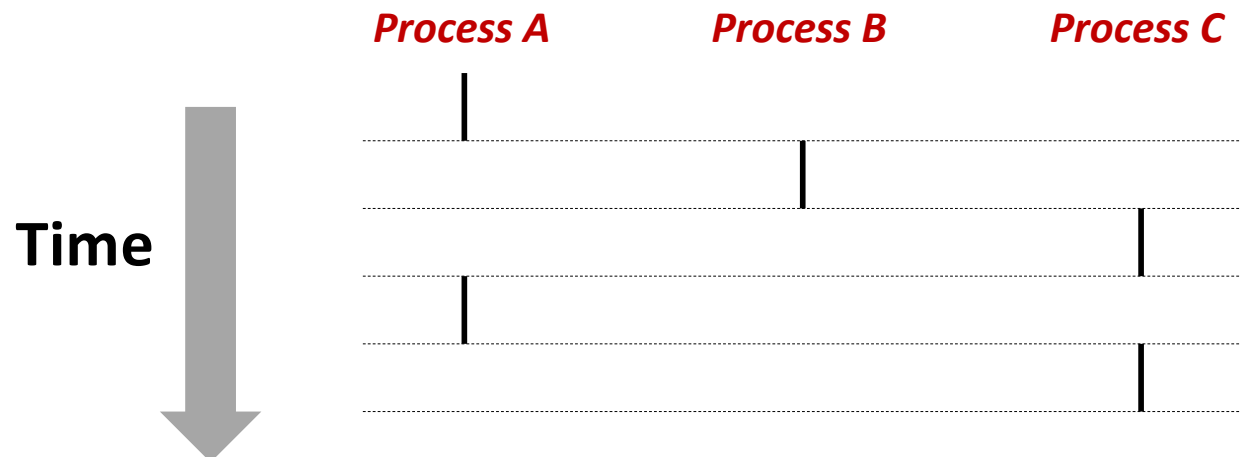
Multiprocessing: The (Modern) Reality



- Multicore processors
 - Multiple CPUs on single chip
 - Share main memory (and some of the caches)
 - Each can execute a separate process
 - Scheduling of processors onto cores done by kernel

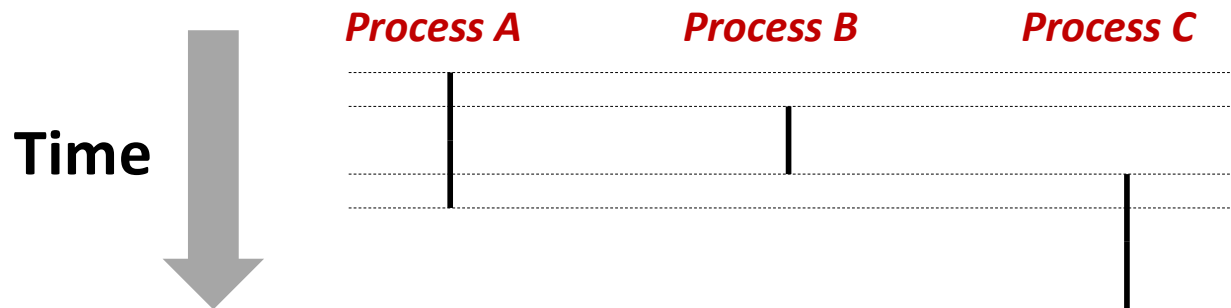
Concurrent Processes

- Each process is a logical control flow
- Two processes run **concurrently** (are concurrent) if their flows overlap in time
- Otherwise, they are **sequential**
- Examples (running on single core):
 - Concurrent: A & B, A & C
 - Sequential: B & C



User View of Concurrent Processes

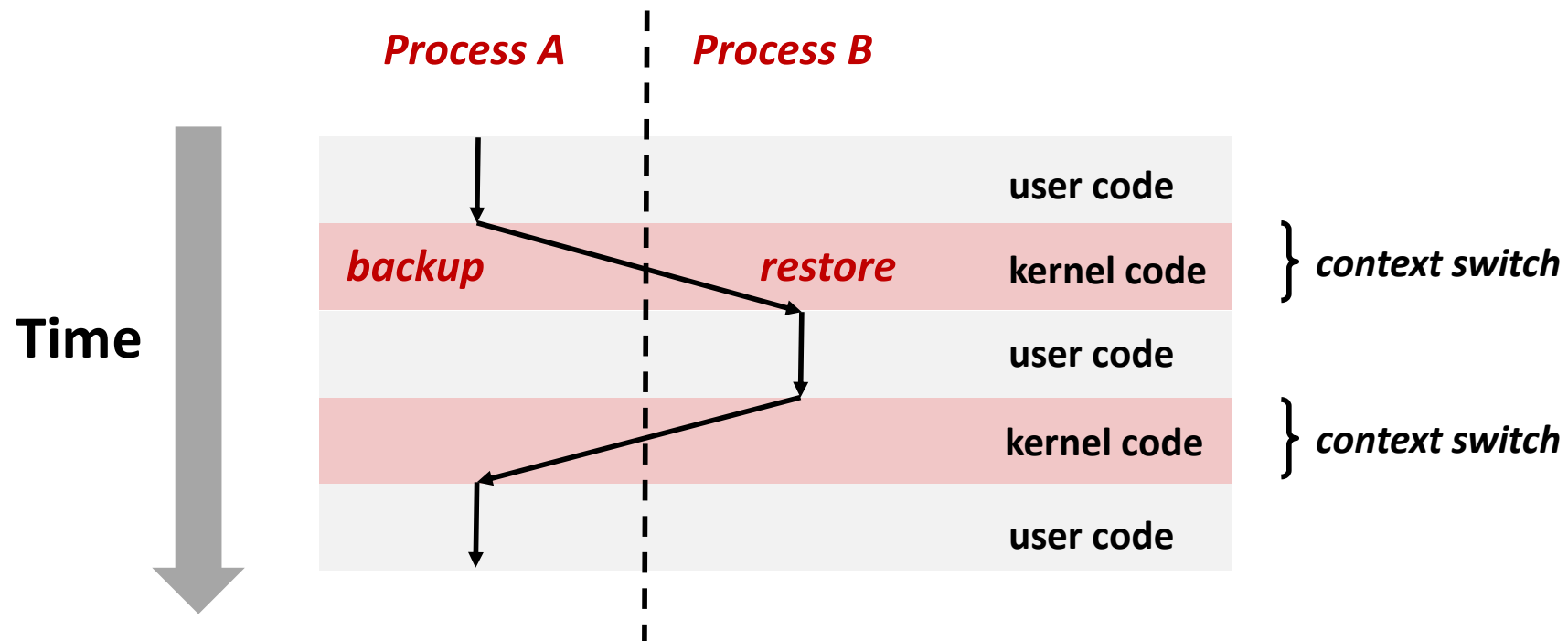
- Control flows for concurrent processes are *physically* disjoint in time
- However, we can think of concurrent processes as *logically* running in parallel with each other



Each process has own logical control flow
but execution time of instructions may vary (why?)

Context Switching

- Processes are managed by a shared chunk of memory-resident OS code called the *kernel*
 - Important: the kernel is not a separate process, but rather runs as part of some existing process
- Control flow passes from one process to another via a *context switch*



Today

- Exceptional Control Flow
- Exceptions
- Processes
- Process Control

System Call Error Handling

- On error, Linux system-level functions typically return -1 and set global variable `errno` to indicate cause
- Hard and fast rule:
 - You must check the return status of every system-level function
 - Only exception is the handful of functions that return `void`
- Example:

```
if ((pid = fork()) < 0) {  
    fprintf(stderr, "fork error: %s\n", strerror(errno));  
    exit(0);  
}
```

Error-reporting functions

- Can simplify somewhat using an error-reporting function:

```
void unix_error(char *msg) /* Unix-style error */
{
    fprintf(stderr, "%s: %s\n", msg, strerror(errno));
    exit(0);
}
```

```
if ((pid = fork()) < 0)
    unix_error("fork error");
```

Error-handling Wrappers

- We simplify the code we present to you even further by using Stevens-style error-handling wrappers:

```
pid_t Fork(void)
{
    pid_t pid;

    if ((pid = fork()) < 0)
        unix_error("Fork error");
    return pid;
}
```

```
pid = Fork();
```

Obtaining Process IDs

- `pid_t getpid(void)`
 - Returns PID of current process
- `pid_t getppid(void)`
 - Returns PID of parent process

Creating and Terminating Processes

- From a programmer's perspective, we can think of a process as being in one of three states
- Running
 - Process is either executing, or waiting to be executed and will eventually be **scheduled** (i.e., chosen to execute) by the kernel
- Stopped
 - Process execution is **suspended** and will not be scheduled until further notice (next lecture when we study signals)
- Terminated
 - Process is stopped permanently

Terminating Processes

- Process becomes terminated for one of three reasons:
 - Returning from the `main` routine
 - Calling the `exit` function
 - Receiving a signal whose default action is to terminate (next lecture)
- `void exit(int status)`
 - Terminates with an *exit* status of `status`
 - Convention: normal return status is 0, nonzero on error
 - Another way to explicitly set the exit status is to return an integer value from the main routine
- `exit` is called **once** but **never** returns.

Creating Processes

- *Parent process* creates a new running child process by calling `fork`
- `int fork(void)`
 - Returns 0 to the child process, child's PID to parent process
 - Child is ***almost*** identical to parent:
 - Child get an identical (but separate) copy of the parent's virtual address space
 - Child gets identical copies of the parent's open file descriptors
 - Child has a different PID than the parent
- `fork` is interesting (and often confusing) because it is called **once** but returns **twice**

fork Example

```
int main()
{
    pid_t pid;
    int x = 1;

    pid = Fork();
    if (pid == 0) { /* Child */
        printf("child : x=%d\n", ++x);
        exit(0);
    }

    /* Parent */
    printf("parent: x=%d\n", --x);
    exit(0);
}
```

fork.c

```
linux> ./fork
parent: x=0
child : x=2
```

- Call once, return twice
- Concurrent execution
 - Can't predict execution order of parent and child
- Duplicate but separate address space
 - `x` has a value of 1 when `fork()` returns in parent and child
 - Subsequent changes to `x` are independent
- Shared open files
 - **stdout** is the same in both parent and child

Modeling fork with Process Graphs

- A **process graph** is a useful tool for capturing the partial ordering of statements in a concurrent program:
 - Each vertex is the execution of a statement
 - $a \rightarrow b$ means a happens before b
 - Edges can be labeled with current value of variables
 - `printf` vertices can be labeled with output
 - Each graph begins with a vertex with no inedges
- Any **topological sort** of the graph corresponds to a feasible total ordering
 - Total ordering of vertices where all edges point from left to right

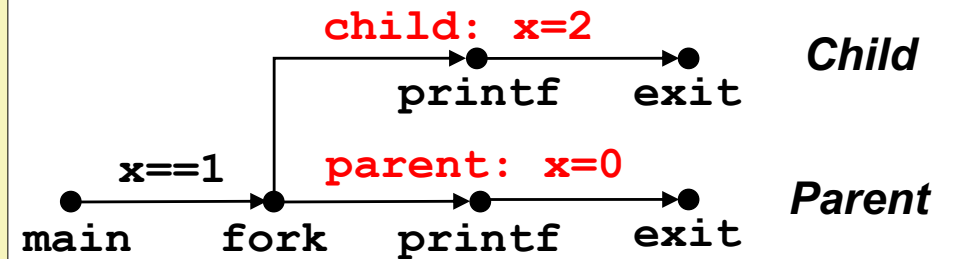
Process Graph Example

```
int main()
{
    pid_t pid;
    int x = 1;

    pid = Fork();
    if (pid == 0) { /* Child */
        printf("child : x=%d\n", ++x);
        exit(0);
    }

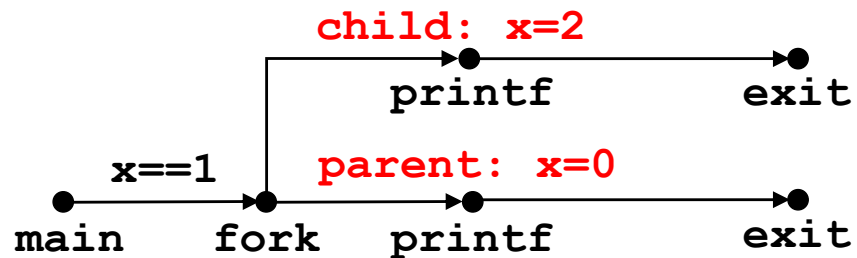
    /* Parent */
    printf("parent: x=%d\n", --x);
    exit(0);
}
```

fork.c

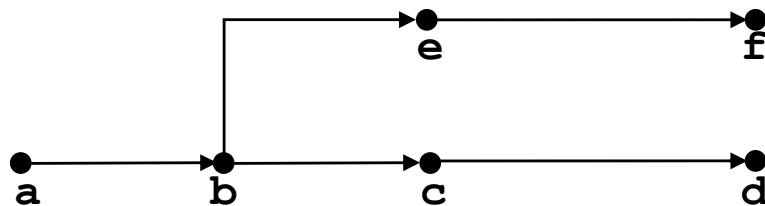


Interpreting Process Graphs

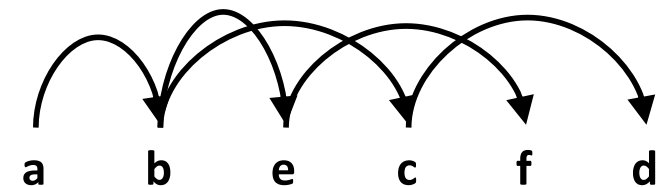
- Original graph:



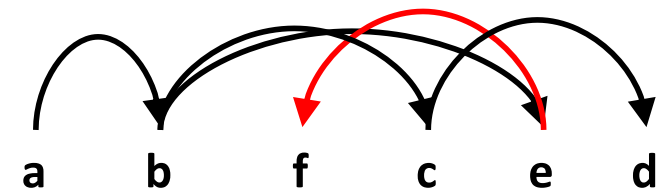
- Relabelled graph:



Feasible total ordering:



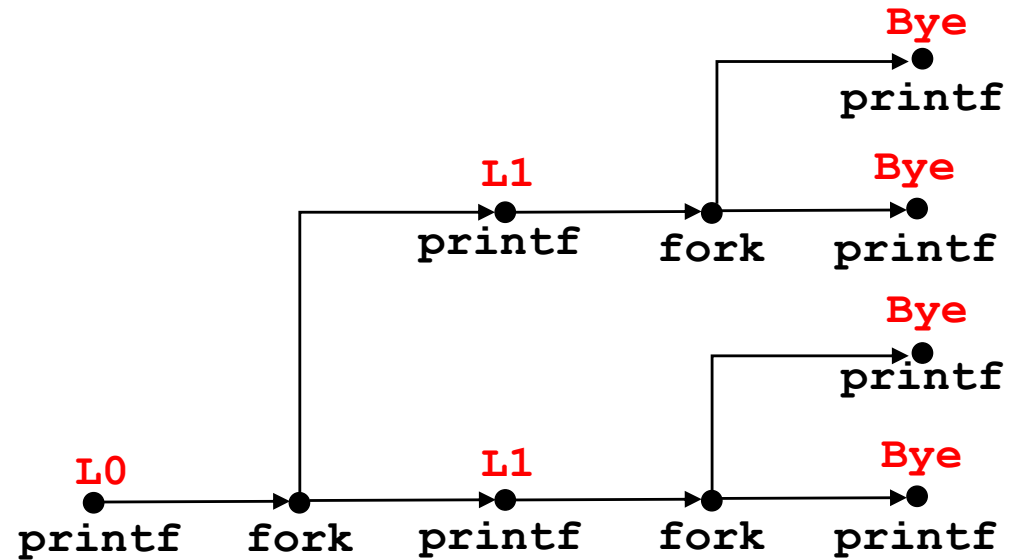
Infeasible total ordering:



fork Example: Two consecutive forks

```
void fork2()
{
    printf("L0\n");
    fork();
    printf("L1\n");
    fork();
    printf("Bye\n");
}
```

forks.c



Feasible output:

L0
L1
Bye
Bye
L1
Bye
Bye

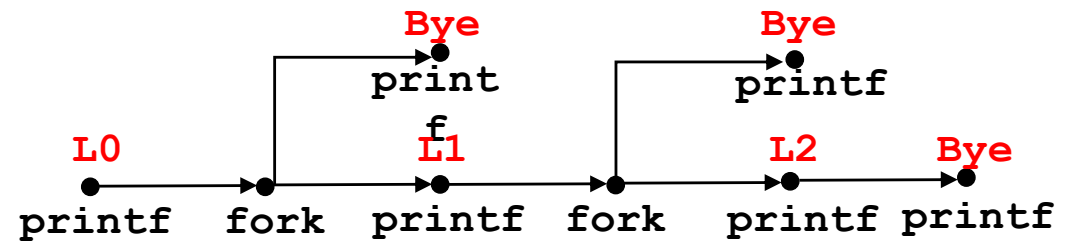
Infeasible output:

L0
Bye
L1
Bye
L1
Bye
Bye

fork Example: Nested forks in parent

```
void fork4()
{
    printf("L0\n");
    if (fork() != 0) {
        printf("L1\n");
        if (fork() != 0) {
            printf("L2\n");
        }
    }
    printf("Bye\n");
}
```

forks.c



Feasible output:

L0
L1
Bye
Bye
L2
Bye

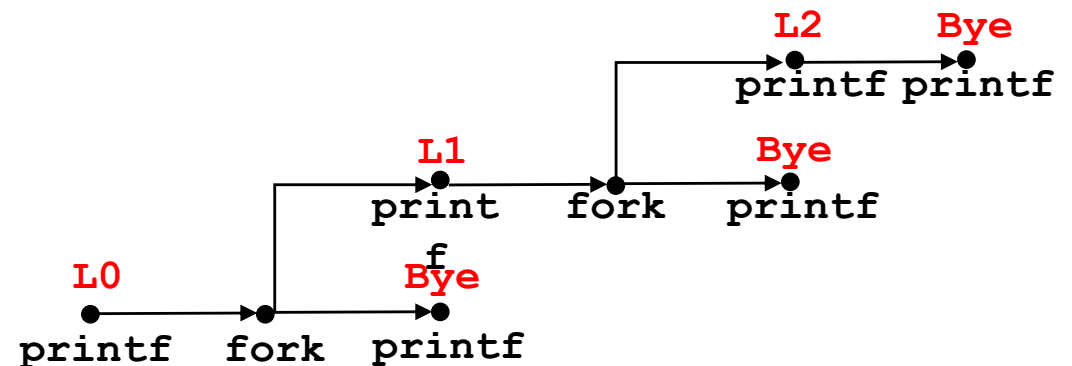
Infeasible output:

L0
Bye
L1
Bye
Bye
L2

fork Example: Nested forks in children

```
void fork5()
{
    printf("L0\n");
    if (fork() == 0) {
        printf("L1\n");
        if (fork() == 0) {
            printf("L2\n");
        }
    }
    printf("Bye\n");
}
```

forks.c



Feasible output:

L0
Bye
L1
L2
Bye
Bye

Infeasible output:

L0
Bye
L1
Bye
Bye
L2

Reaping Child Processes

- Idea
 - When process terminates, it still consumes system resources
 - Examples: Exit status, various OS tables
 - Called a “zombie”
 - Living corpse, half alive and half dead
- Reaping
 - Performed by parent on terminated child (using `wait` or `waitpid`)
 - Parent is given exit status information
 - Kernel then deletes zombie child process
- What if parent doesn't reap?
 - If any parent terminates without reaping a child, then the orphaned child will be reaped by `init` process (`pid == 1`)
 - So, only need explicit reaping in long-running processes
 - e.g., shells and servers

Zombie Example

```
void fork7() {
    if (fork() == 0) {
        /* Child */
        printf("Terminating Child, PID = %d\n", getpid());
        exit(0);
    } else {
        printf("Running Parent, PID = %d\n", getpid());
        while (1); /* Infinite loop */
    }
}
```

forks.c

```
linux> ./forks 7 &
[1] 6639
Running Parent, PID = 6639
Terminating Child, PID = 6640
linux> ps
  PID TTY          TIME CMD
 6585 tttyp9        00:00:00 tcsh
 6639 tttyp9        00:00:03 forks
 6640 tttyp9        00:00:00 forks <defunct>
 6641 tttyp9        00:00:00 ps
linux> kill 6639
[1]      Terminated
linux> ps
  PID TTY          TIME CMD
 6585 tttyp9        00:00:00 tcsh
 6642 tttyp9        00:00:00 ps
```

ps shows child process as “defunct”
(i.e., a zombie)

Killing parent allows child to be
reaped by `init`

Non-terminating Child Example

```
void fork8() {
    if (fork() == 0) {
        /* Child */
        printf("Running Child, PID = %d\n",
            getpid());
        while (1); /* Infinite loop */
    } else {
        printf("Terminating Parent, PID = %d\n",
            getpid());
        exit(0);
    }
}
```

forks.c

```
linux> ./forks 8
Terminating Parent, PID = 6675
Running Child, PID = 6676
```

```
linux> ps
  PID TTY          TIME CMD
 6585 tttyp9      00:00:00 tcsh
 6676 tttyp9      00:00:06 forks
 6677 tttyp9      00:00:00 ps
```

```
linux> kill 6676
```

```
linux> ps
  PID TTY          TIME CMD
 6585 tttyp9      00:00:00 tcsh
 6678 tttyp9      00:00:00 ps
```

Child process still active even though parent has terminated

Must kill child explicitly, or else will keep running indefinitely

wait: Synchronizing with Children

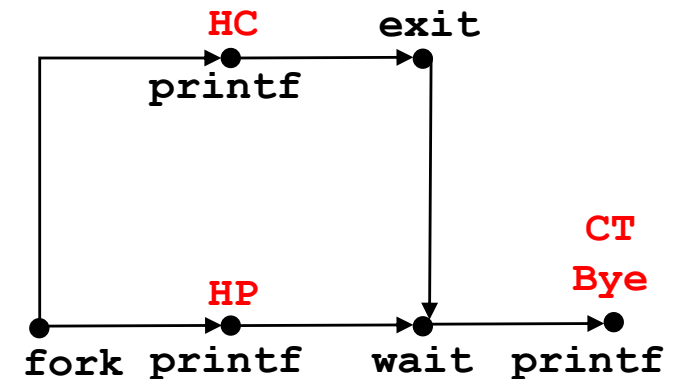
- Parent reaps a child by calling the `wait` function
- `pid_t wait(int* child_status)`
 - Suspends current process until one of its children terminates
 - Return value is the **pid** of the child process that terminated
 - If **child_status != NULL**, then the integer it points to will be set to a value that indicates reason the child terminated and the exit status:
 - Checked using macros defined in `wait.h`
 - `WIFEXITED`, `WEXITSTATUS`, `WIFSIGNALED`, `WTERMSIG`, `WIFSTOPPED`, `WSTOPSIG`, `WIFCONTINUED`
 - See textbook for details

wait: Synchronizing with Children

```
void fork9() {
    int child_status;

    if (fork() == 0) {
        printf("HC: hello from child\n");
        exit(0);
    } else {
        printf("HP: hello from parent\n");
        wait(&child_status);
        printf("CT: child has terminated\n");
    }
    printf("Bye\n");
}
```

forks.c



Feasible output:

HC
HP
CT
Bye

Infeasible output:

HP
CT
Bye
HC

Another wait Example

- If multiple children completed, will take in arbitrary order
- Can use macros `WIFEXITED` and `WEXITSTATUS` to get information about exit status

```
void fork10() {
    pid_t pid[N];
    int i, child_status;
    for (i = 0; i < N; i++)
        if ((pid[i] = fork()) == 0) {
            exit(100+i); /* Child */
        }
    for (i = 0; i < N; i++) { /* Parent */
        pid_t wpid = wait(&child_status);
        if (WIFEXITED(child_status))
            printf("Child %d terminated with exit status %d\n",
                wpid, WEXITSTATUS(child_status));
        else
            printf("Child %d terminate abnormally\n", wpid);
    }
}
```

forks.c

waitpid: Waiting for a Specific Process

- `pid_t waitpid(pid_t pid, int* status, int options)`
 - Suspends current process until specific process terminates
 - Various options (see textbook)

```
void fork11() {
    pid_t pid[N];
    int i;
    int child_status;

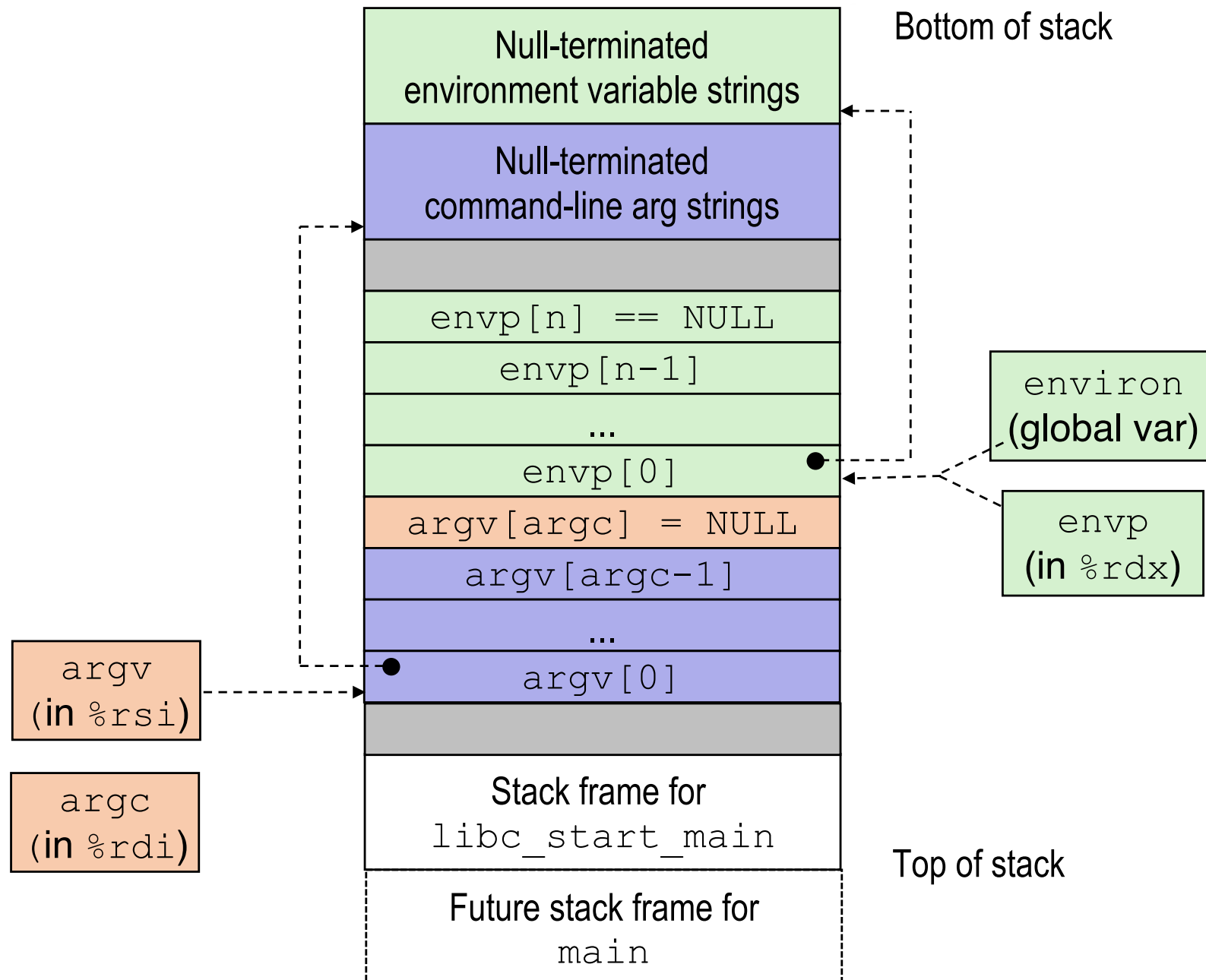
    for (i = 0; i < N; i++)
        if ((pid[i] = fork()) == 0)
            exit(100+i); /* Child */
    for (i = N-1; i >= 0; i--) {
        pid_t wpid = waitpid(pid[i], &child_status, 0);
        if (WIFEXITED(child_status))
            printf("Child %d terminated with exit status %d\n",
                wpid, WEXITSTATUS(child_status));
        else
            printf("Child %d terminate abnormally\n", wpid);
    }
}
```

forks.c

execve: Loading and Running Programs

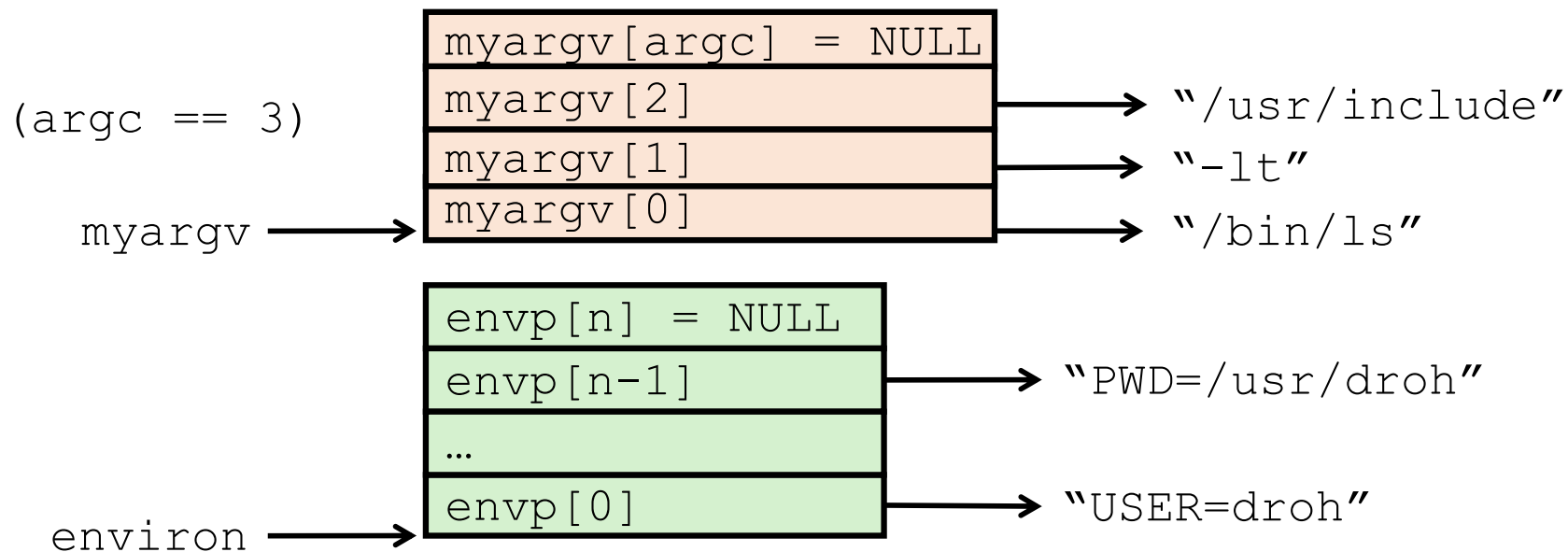
- `int execve(char *filename, char *argv[], char *envp[])`
- Loads and runs in the current process:
 - Executable file **filename**
 - Can be object file or script file beginning with `#!/interpreter` (e.g., `#!/bin/bash`)
 - ...with argument list **argv**
 - By convention **argv[0]==filename**
 - ...and environment variable list **envp**
 - “name=value” strings (e.g., `USER=droh`)
 - `getenv`, `putenv`, `putenv`
- Overwrites code, data, and stack
 - Retains PID, open files and signal context
- Called **once** and **never** returns
 - ...except if there is an error

Structure of Stack When a Program Starts



execve Example

- Executes `/bin/ls -lt /usr/include` in child process using current environment



```

if ((pid = Fork()) == 0) {    /* Child runs program */
    if (execve(myargv[0], myargv, environ) < 0) {
        printf("%s: Command not found.\n", myargv[0]);
        exit(1);
    }
}

```

Summary

- Exceptions
 - Events that require nonstandard control flow
 - Generated externally (interrupts) or internally (traps and faults)
- Processes
 - At any given time, system has multiple active processes
 - Only one can execute at a time on a single core, though
 - Each process appears to have total control of processor + private memory space

Summary (cont.)

- Spawning processes
 - Call `fork`
 - One call, two returns
- Process completion
 - Call `exit`
 - One call, no return
- Reaping and waiting for processes
 - Call `wait` or `waitpid`
- Loading and running programs
 - Call `execve` (or variant)
 - One call, (normally) no return