

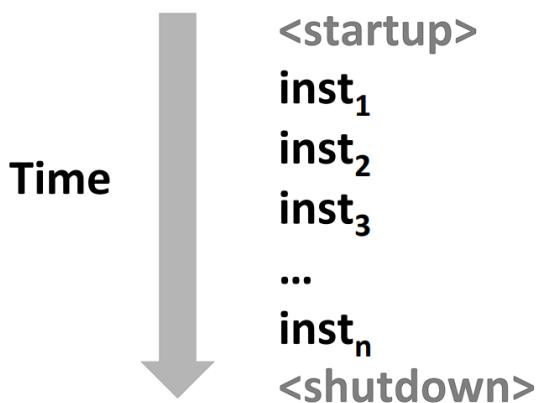
Exceptional Control Flow

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

- Jump and branches
- Call and return

Change **program state**.

But 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

So, system needs mechanisms for **ECF**

ECF

- **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
 - 4. **Nonlocal jumps**: `setjmp()` and `longjmp()`
 - Implemented by C runtime library

that exist at all levels of a computer system. We start with exceptions, which lie at the intersection of the hardware and the operating system. We also discuss system calls, which are exceptions that provide applications with entry points into the operating system. We then move up a level of abstraction and describe processes and signals, which lie at the intersection of applications and the operating system. Finally, we discuss nonlocal jumps, which are an application-level form of ECF.

Exceptions

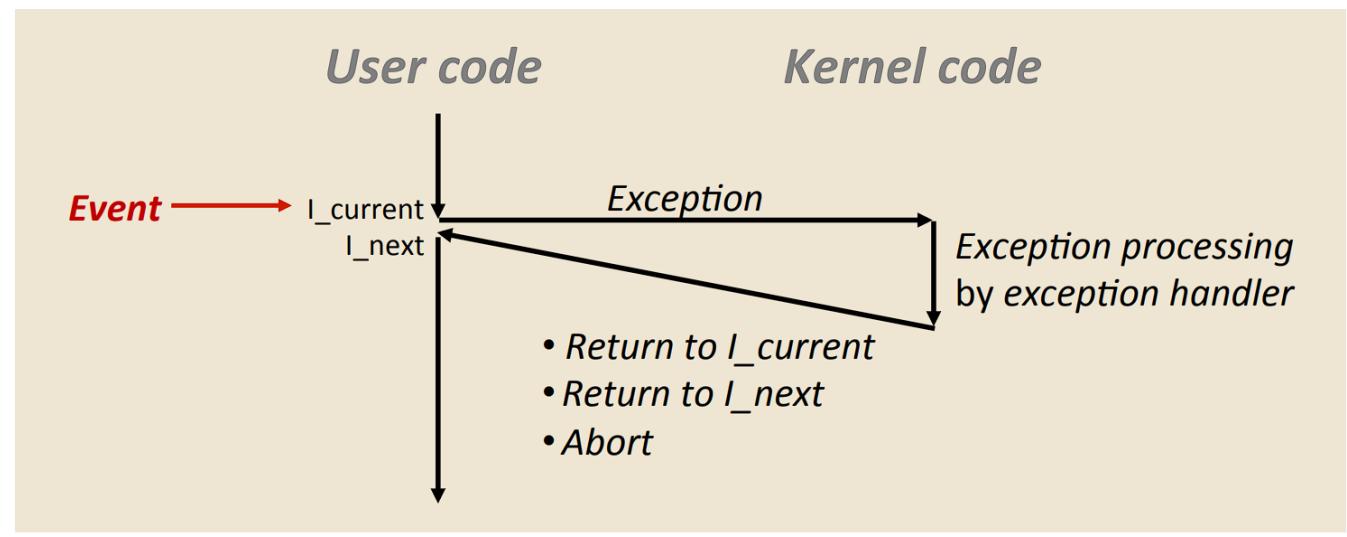
An exception is a **transfer of control to the OS kernel** in response to some **event** (i.e., change in processor state). An exception is an abrupt change in the control flow in response to some change in the **processor's state**.

Can be divided into four classes:

interrupts,traps,faults,aborts.

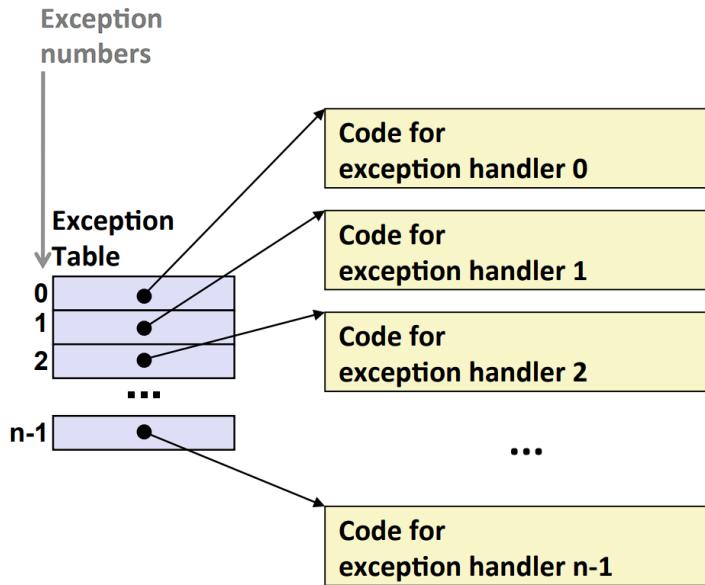
Each type of possible exception in a system is assigned a unique nonnegative integer *exception number*. Some of these numbers are assigned by the designers of the processor. Other numbers are assigned by the designers of the operating system *kernel* (the memory-resident part of the operating system). Examples of the former include divide by zero, page faults, memory access violations, breakpoints, and arithmetic overflows. Examples of the latter include system calls and signals from external I/O devices.

- 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



When control is being transferred from a user program to the kernel, all of these items are pushed onto the kernel's stack rather than onto the user's stack.

Exception Tables



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

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

- **Interrupt** – caused by an event *external* to current running program
 - E.g., key press, disk I/O
 - Asynchronous to current program
 - Can handle interrupt on any convenient instruction
 - "Whenever it's convenient, just don't wait too long"

Synchronous Exceptions

- **Exception** – caused by some event *during* execution of one instruction of current running program
 - E.g., memory error, bus error, illegal instruction, raised exception
 - Synchronous
 - Must handle exception *precisely* on instruction that causes exception
 - “Drop whatever you are doing and act now”
- **Trap** – action of servicing interrupt or exception by hardware jump to “interrupt or trap handler” code

■ **Caused by events that occur as a result of executing an instruction:**

- **Traps**
 - Intentional
 - 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

Syscall

System Calls

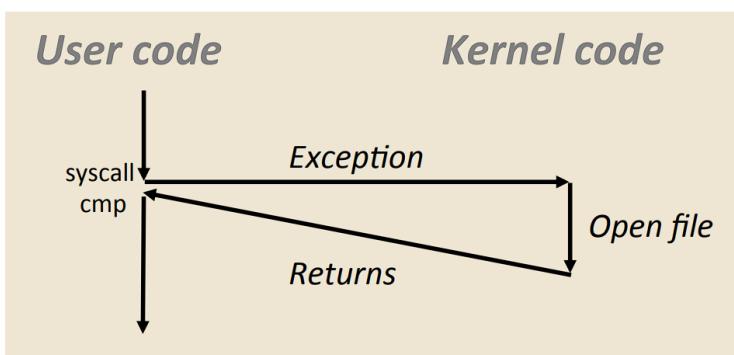
- Each x86-64 system call has a unique ID number
- Examples:

Number	Name	Description
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` function, which invokes system call instruction `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    $0xfffffffffffff001,%rax  
...  
e5dfa: c3                   retq
```



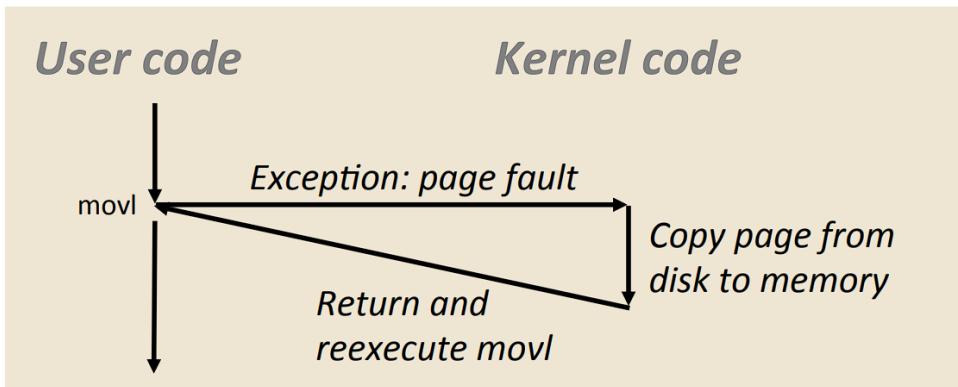
- `%rax` contains syscall number
- Other arguments in `%rdi`, `%rsi`, `%rdx`, `%r10`, `%r8`, `%r9`
- Return value in `%rax`
- Negative value is an error corresponding to negative `errno`

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
```



Processes

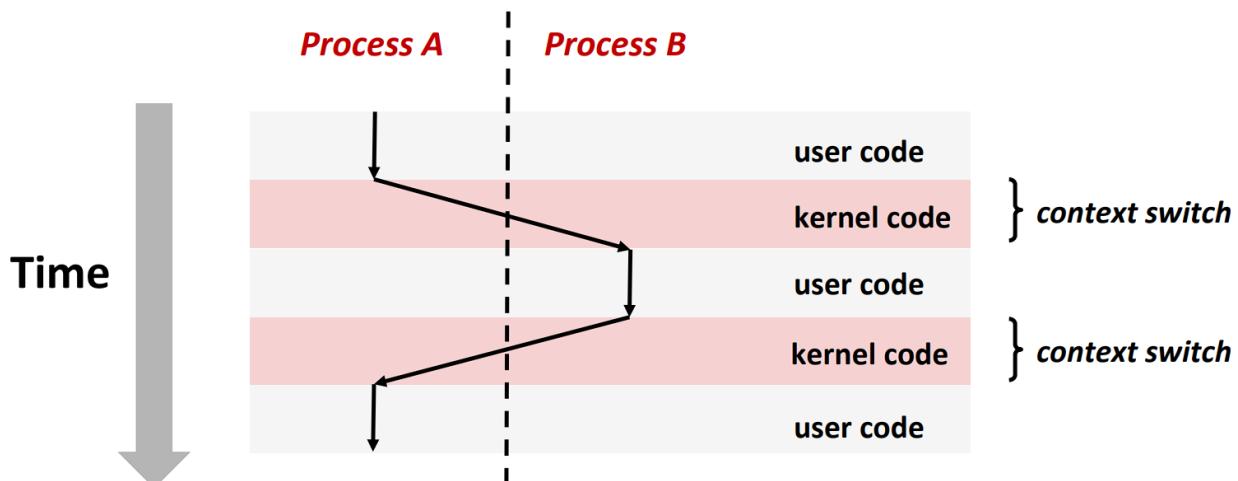
errno: a global int which declares what type of error has happened.

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*

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*



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`

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
-
- **Process becomes terminated for one of three reasons:**
 - Receiving a signal whose default action is to terminate (next lecture)
 - Returning from the `main` routine
 - Calling the `exit` function

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 ttys000  00:00:00 tcsh
 6639 ttys000  00:00:03 forks
 6640 ttys000  00:00:00 forks <defunct>
 6641 ttys000  00:00:00 ps
```

```
linux> kill 6639
[1]  Terminated
```

```
linux> ps
  PID TTY      TIME CMD
 6585 ttys000  00:00:00 tcsh
 6642 ttys000  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

- `int 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

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`, `printenv`

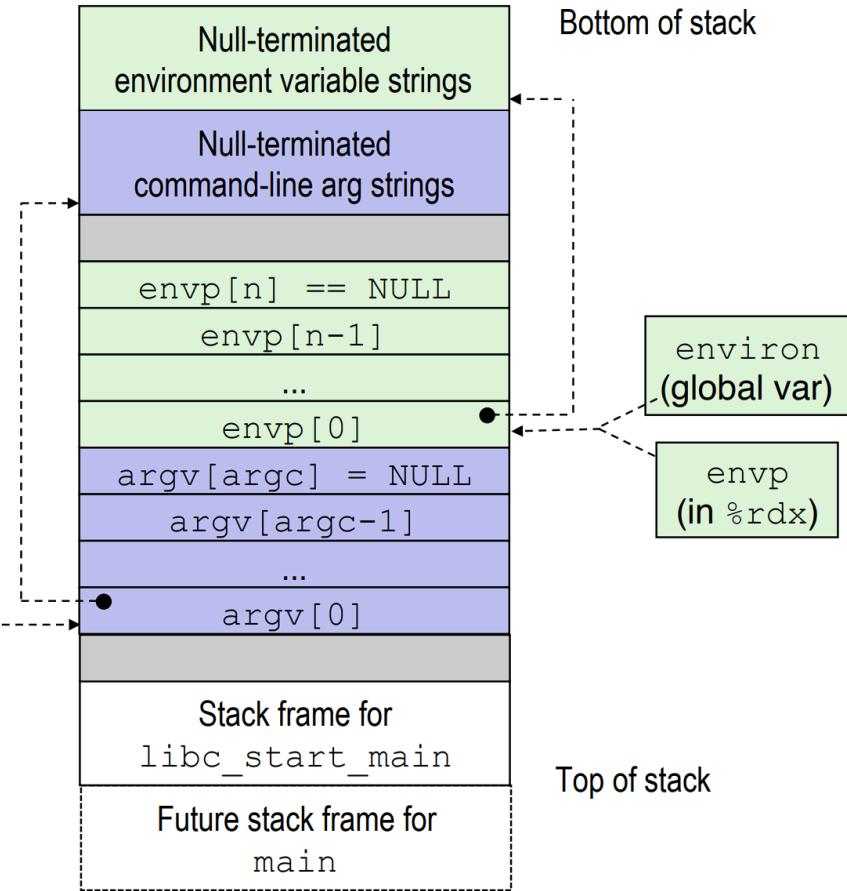
- Overwrites code, data, and stack

- Retains PID, open files and signal context

- Called once and never returns

- ...except if there is an error

Structure of the stack when a new program starts



Signal

■ But what about background jobs?

- Will become zombies when they terminate
- Will never be reaped because shell (typically) will not terminate
- Will create a memory leak that could run the kernel out of memory

Background process will not be terminated after running, **kernel** need to send signals to process which will be terminated.

A signal is a small message which **notifies** a process that

an **event** of some type have occurred in the system.

- Akin to exceptions and interrupts
- Sent from the kernel (sometimes at the request of another process) to a process
- Signal type is identified by small integer ID's (1-30)
- Only information in a signal is its ID and the fact that it arrived

ID	Name	Default Action	Corresponding Event
2	SIGINT	Terminate	User typed ctrl-c
9	SIGKILL	Terminate	Kill program (cannot override or ignore)
11	SIGSEGV	Terminate & Dump	Segmentation violation
14	SIGALRM	Terminate	Timer signal
17	SIGCHLD	Ignore	Child stopped or terminated

Sending signal

Kernel sends a signal to a destination process by **updating** some **state** in the **context** of the destination process.

Why sending signal

Kernel sends a signal for one of the following reasons:

- Kernel has detected a system event such as divide-by-zero (SIGFPE) or the termination of a child process (SIGCHLD)
- Another process has invoked the `kill` system call to explicitly request the kernel to send a signal to the destination process

Sending Signals with /bin/kill Program

/bin/kill program sends arbitrary signal to a process or process group

Examples

- `/bin/kill -9 24818`
Send SIGKILL to process 24818
- `/bin/kill -9 -24817`
Send SIGKILL to every process in process group 24817

```
linux> ./forks 16
Child1: pid=24818 pgrp=24817
Child2: pid=24819 pgrp=24817

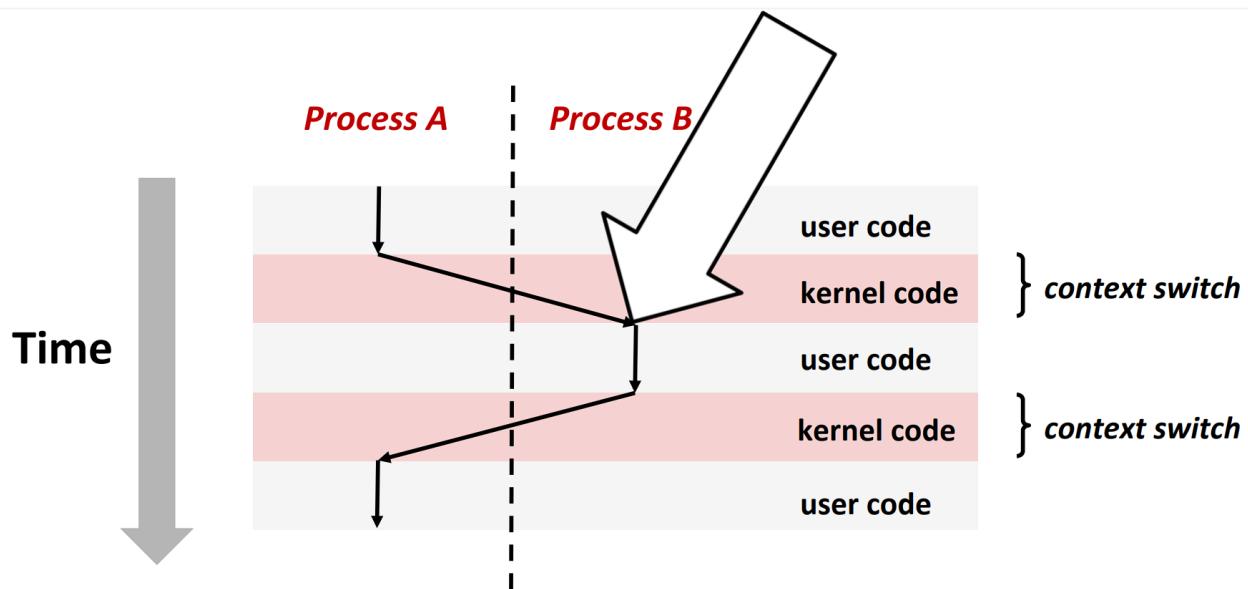
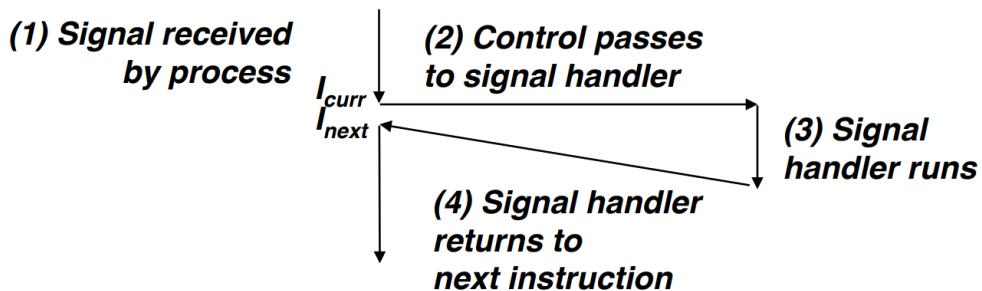
linux> ps
 PID TTY          TIME CMD
24788 pts/2        00:00:00 tcsh
24818 pts/2        00:00:02 forks
24819 pts/2        00:00:02 forks
24820 pts/2        00:00:00 ps
linux> /bin/kill -9 -24817
linux> ps
 PID TTY          TIME CMD
24788 pts/2        00:00:00 tcsh
24823 pts/2        00:00:00 ps
linux>
```

Sending signals(SIGINT/SIGTSTP) from the Keyboard will only act on the foreground process group.

Receiving a Signal

Signal Concepts: Receiving a Signal

- A destination process **receives** a signal when it is forced by the kernel to react in some way to the delivery of the signal
- Some possible ways to react:
 - **Ignore** the signal (do nothing)
 - **Terminate** the process (with optional core dump)
 - **Catch** the signal by executing a user-level function called **signal handler**
 - Akin to a hardware exception handler being called in response to an asynchronous interrupt:



Important: All context switches are initiated by calling some exception handler.

Receiving Signals

- Suppose kernel is returning from an exception handler and is ready to pass control to process p
- Kernel computes $\text{pnb} = \text{pending} \& \sim\text{blocked}$
 - The set of pending nonblocked signals for process p
- If ($\text{pnb} == 0$)
 - Pass control to next instruction in the logical flow for p
- Else
 - Choose least nonzero bit k in pnb and force process p to **receive** signal k
 - The receipt of the signal triggers some **action** by p
 - Repeat for all nonzero k in pnb
 - Pass control to next instruction in logical flow for p

当内核把进程 p 从内核模式切换到用户模式时(例如, 从系统调用返回或是完成了一次上下文切换), 它会检查进程 p 的未被阻塞的待处理信号的集合($\text{pending} \& \sim\text{blocked}$)。

Pending(挂起) and Blocked Signals

Pending: sent but not yet received.

- There can be at most one pending signal of any particular type
- Important: Signals are not queued
 - If a process has a pending signal of type k , then subsequent signals of type k that are sent to that process are discarded

- A process can **block** the receipt of certain signals

- Blocked signals can be delivered, but will not be received until the signal is unblocked

A pending signal is received at most once.

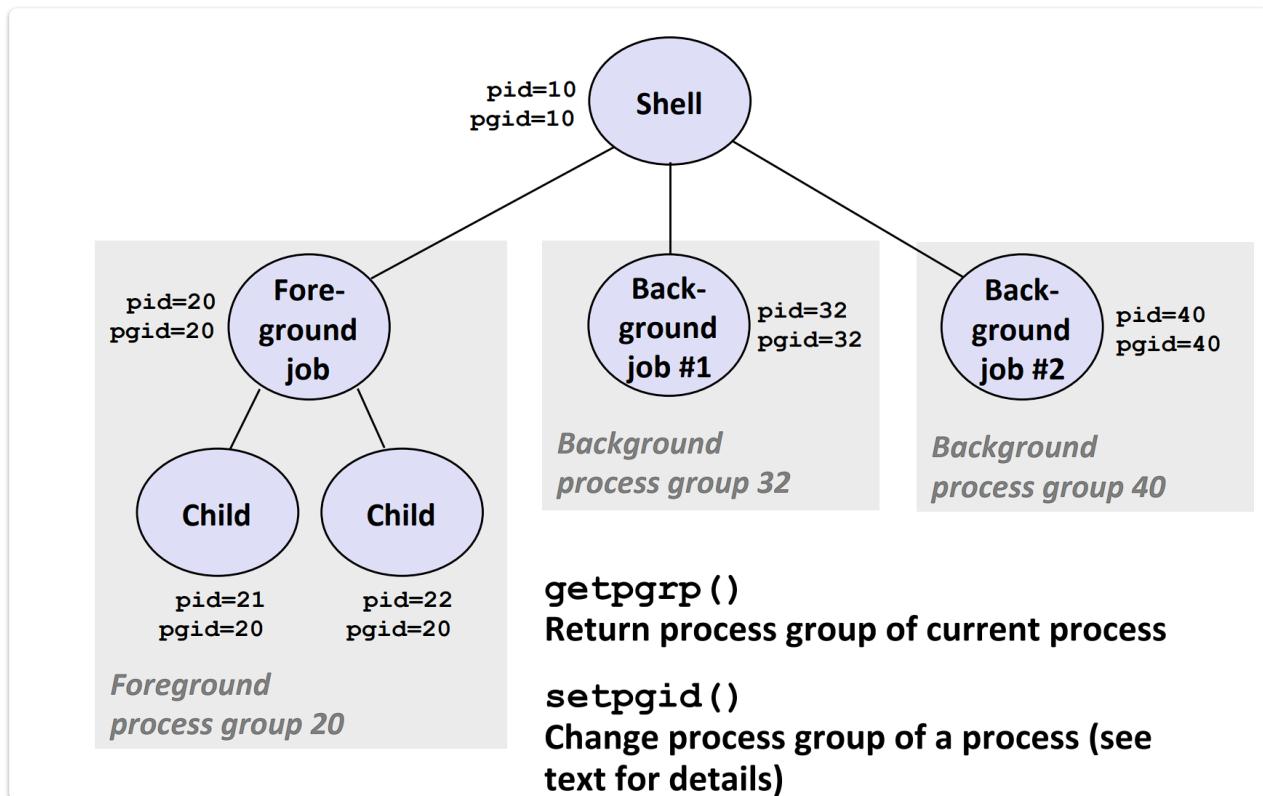
Pending / Blocked Bits

Signal Concepts: Pending/Blocked Bits

- Kernel maintains **pending** and **blocked** bit vectors in the context of each process

- **pending**: represents the set of pending signals
 - Kernel sets bit k in **pending** when a signal of type k is delivered
 - Kernel clears bit k in **pending** when a signal of type k is received
- **blocked**: represents the set of blocked signals
 - Can be set and cleared by using the **sigprocmask** function
 - Also referred to as the *signal mask*.

- Every process belongs to exactly one process group.



Default Actions of signal

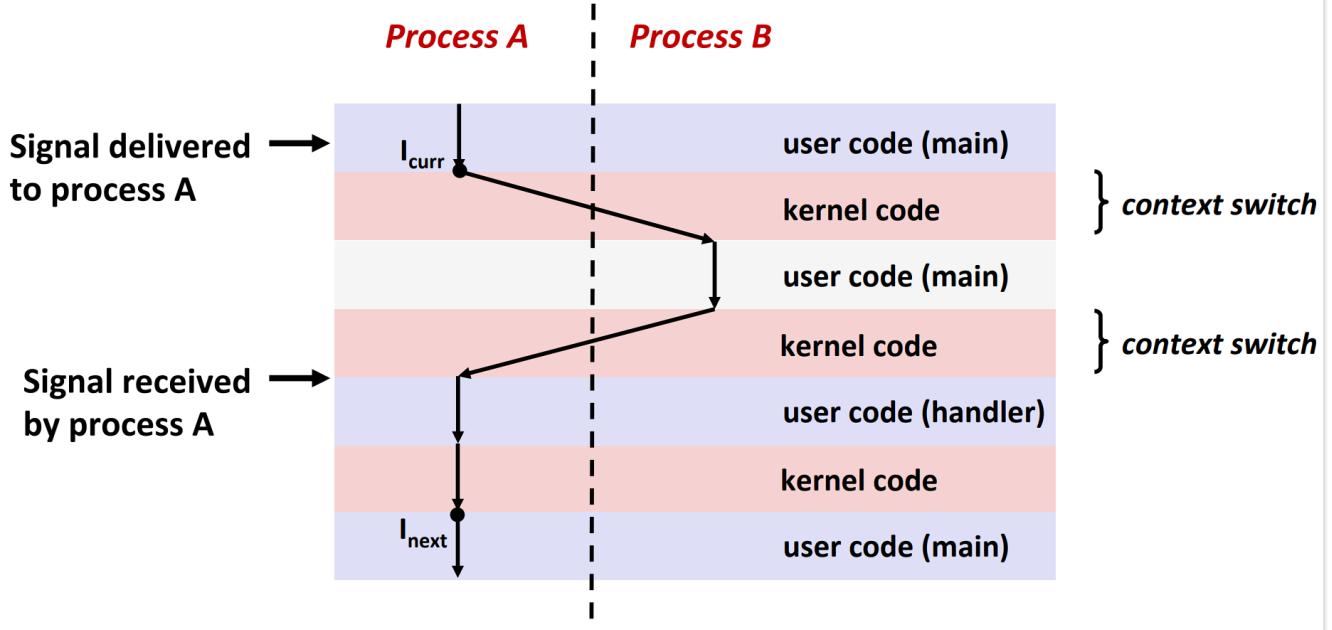
Each signal type has a predefined **default action**, which is one of:

- The process terminates.
- The process terminates and dumps core(take a snap shot of the memory).
- The process stops until restarted by a SIGCONT signal.
- The process ignores the signal.

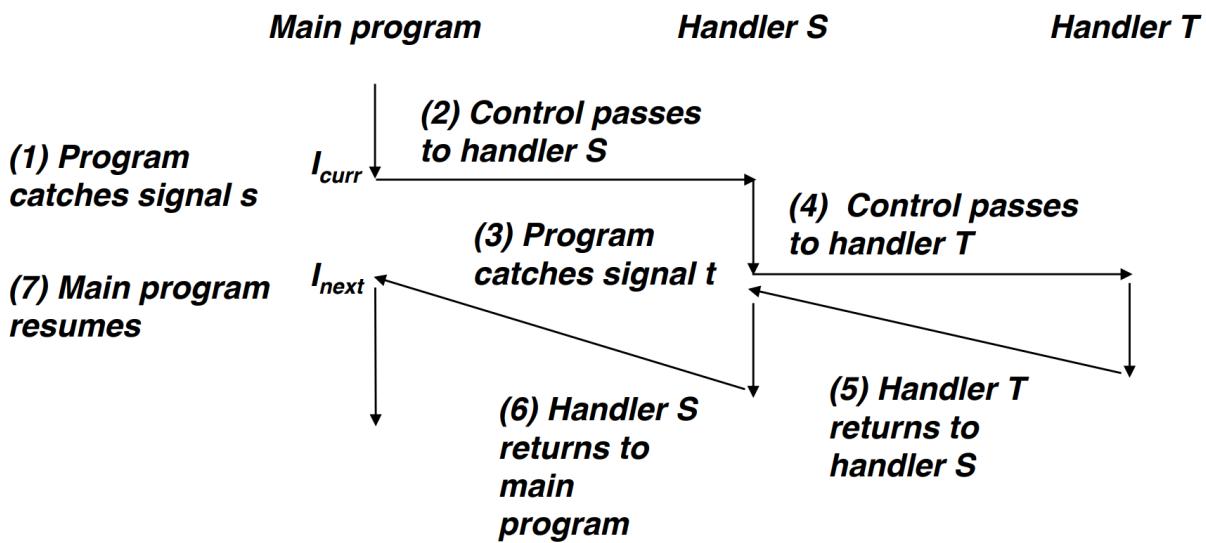
Install Signal Handlers

- The **signal** function modifies the default action associated with the receipt of signal **signum**:
 - `handler_t *signal(int signum, handler_t *handler)`
- Different values for **handler**:
 - **SIG_IGN**: ignore signals of type **signum**
 - **SIG_DFL**: revert to the default action on receipt of signals of type **signum**
 - Otherwise, **handler** is the address of a user-level **signal handler**
 - Called when process receives signal of type **signum**
 - Referred to as “**installing**” the handler
 - Executing handler is called “**catching**” or “**handling**” the signal
 - When the handler executes its return statement, control passes back to instruction in the control flow of the process that was interrupted by receipt of the signal

A signal handler is a separate logical flow(not process) that runs concurrently with the main program.



■ Handlers can be interrupted by other handlers



Blocking and unblocking signals

Implicit

- Kernel blocks any pending signals of type currently being handled.

Blocking and Unblocking Signals

- **Implicit blocking mechanism**
 - Kernel blocks any pending signals of type currently being handled.
 - E.g., A SIGINT handler can't be interrupted by another SIGINT
- **Explicit blocking and unblocking mechanism**
 - `sigprocmask` function
- **Supporting functions**
 - `sigemptyset` – Create empty set
 - `sigfillset` – Add every signal number to set
 - `sigaddset` – Add signal number to set
 - `sigdelset` – Delete signal number from set

Safe Signal Handling

- **Handlers are tricky because they are concurrent with main program and share the same global data structures.**
 - Shared data structures can become corrupted.
- **We'll explore concurrency issues later in the term.**
- **For now here are some guidelines to help you avoid trouble.**