# Exceptional Control Flow: Exceptions and Processes
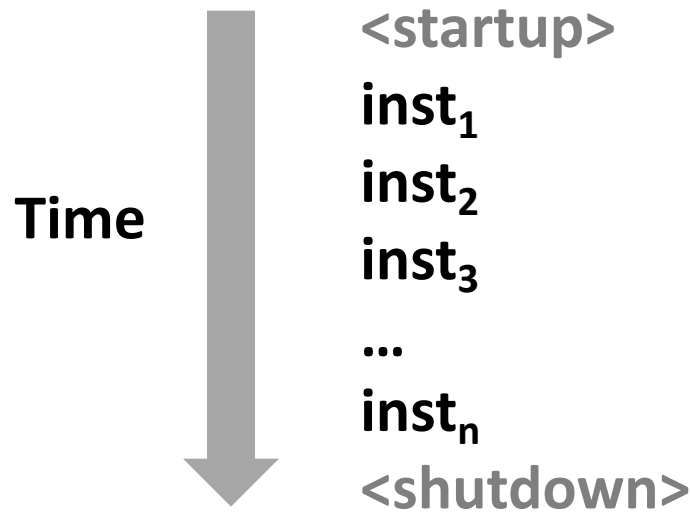
# Today

- **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*

Time ⟶

<startup>
inst$_1$
inst$_2$
inst$_3$
…
inst$_n$

# 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 "exceptional control flow"**
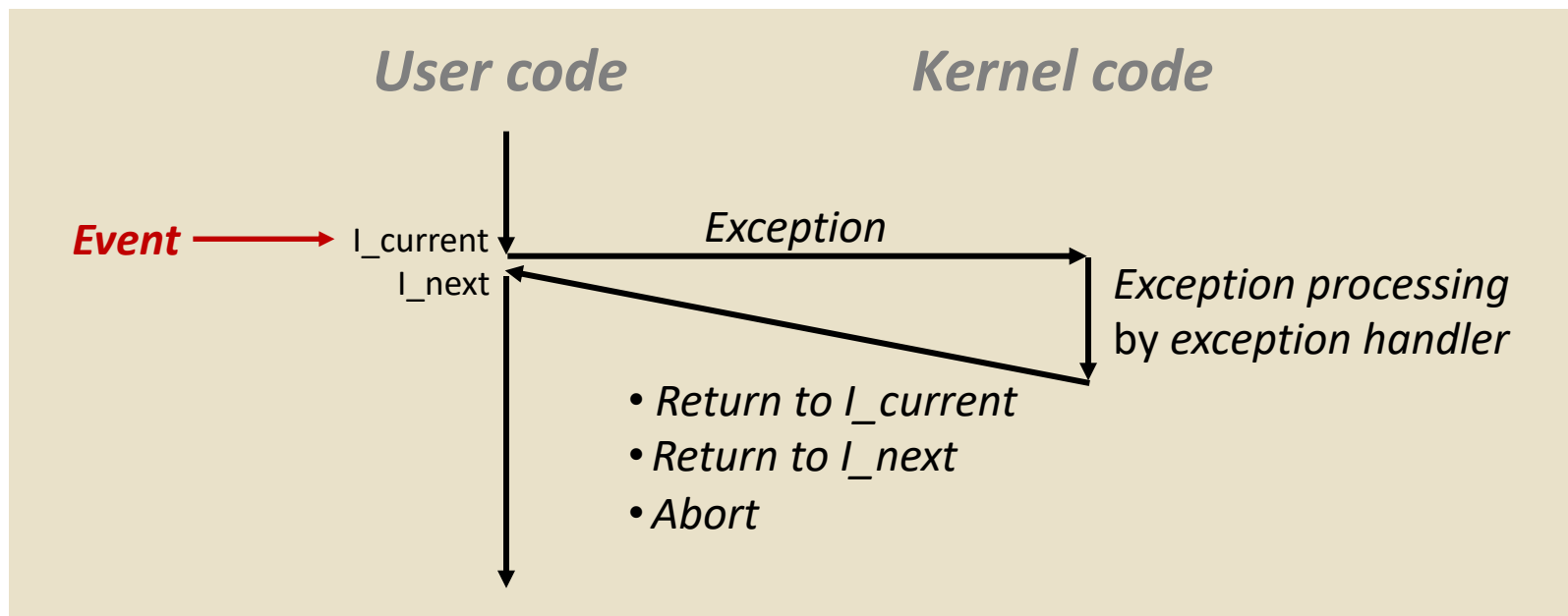
# Exceptional Control Flow

- **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

  - 4. **Nonlocal jumps**: `setjmp()` and `longjmp()`

    - cross-function jumps ,  goto only happen within a function
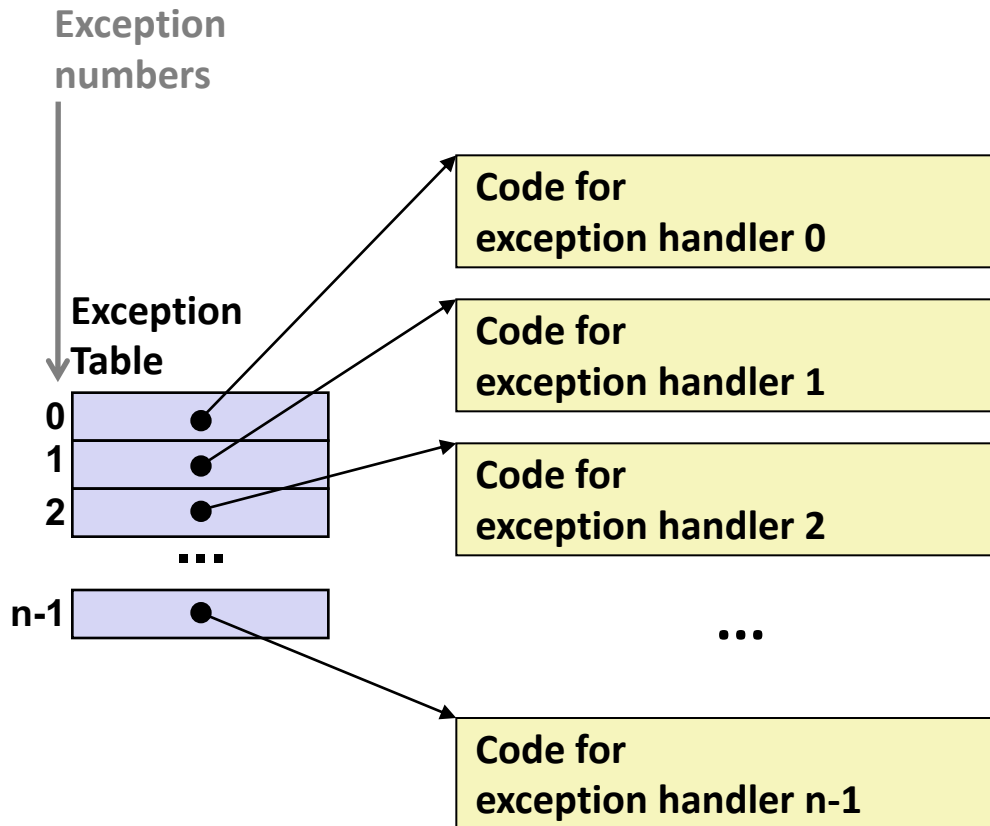    - Implemented by C runtime library

# 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
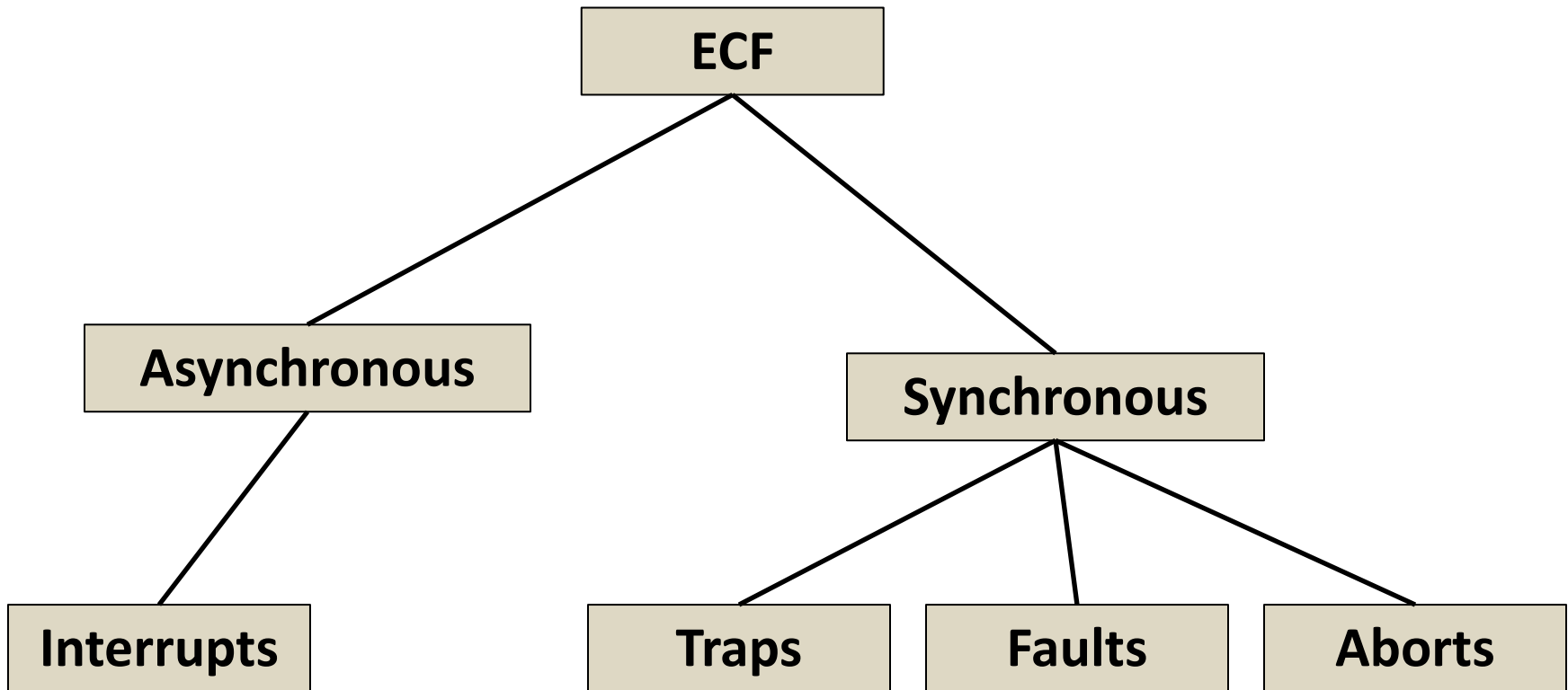
*User code*          *Kernel code*

*Event* ———→ I_current ————— *Exception* ————————→
            I_next ←⟍                          *Exception processing*
                     ⟍                          by *exception handler*
                      ⟍—————————————→

• *Return to I_current*
• *Return to I_next*
• *Abort*

# Exception Tables

**Exception numbers**

**Exception Table**

0
1
2

...

n-1

Code for
exception handler 0

Code for
exception handler 1

Code for
exception handler 2

...

Code for
exception handler n-1

- **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**

# (partial) Taxonomy

# 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
  - 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
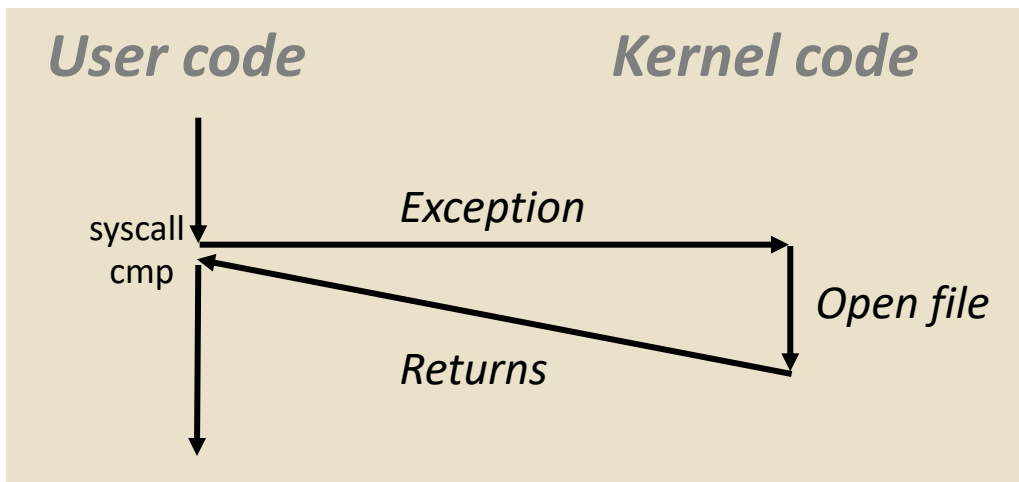  - Aborts current program

# System Calls

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

| 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
```



*User code*     *Kernel code*

syscall → *Exception* → *Open file*
cmp ← *Returns*

- `%eax` 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`

# System Call Example: Opening File

- User calls: `open(filename, options)`
- Calls __o...

```
000000000
...
e5d79:                                         l #2
e5d7e:                                         %rax
e5d80:
...
e5dfa:
```
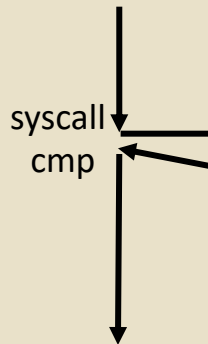
*User code*

syscall
cmp

A box overlays the slide:

**Almost like a function call**
- **Transfer of control**
- **On return, executes next instruction**
- **Passes arguments using calling convention**
- **Gets result in `%rax`**

**One Important exception!**
- **Executed by Kernel**
- **Different set of privileges**
- **And other differences:**
  - **E.g., "address" of "function" is in `%rax`**
  - **Uses `errno`**
  - **Etc.**

number
rdi,
r8,%r9

ror
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;
}
```
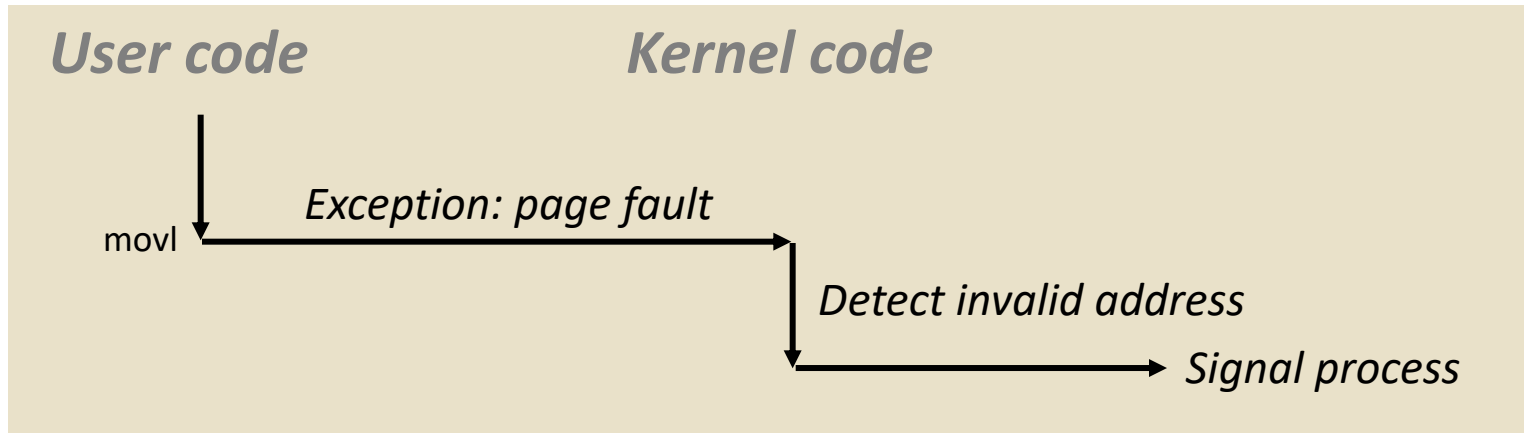
```
80483b7:        c7 05 10 9d 04 08 0d   movl    $0xd,0x8049d10
```

**User code**          **Kernel code**

movl

Exception: page fault

Copy page from disk to memory

Return and reexecute movl

# Fault Example: Invalid Memory Reference

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

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

**User code**                    **Kernel code**

movl

Exception: page fault

Detect invalid address

Signal process

- Sends **SIGSEGV** signal to user process
- User process exits with "segmentation fault"

# Today

- **Exceptional Control Flow**

- **Exceptions**

- **Processes**

- **Process Control**

# Processes

- **Definition: A *process* is 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*

**Memory**

| Stack |
| Heap |
| Data |
| Code |

**CPU**

| Registers |

# Multiprocessing: The Illusion

| Memory | | Memory | | | Memory |
|---|---|---|---|---|---|
| Stack | | Stack | | | Stack |
| Heap | | Heap | | ••• | Heap |
| Data | | Data | | | Data |
| Code | | Code | | | Code |
| **CPU** | | **CPU** | | | **CPU** |
| Registers | | Registers | | | Registers |

- **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 389 processes, 2 of which are active

  ▪ Identified by Process ID (PID)

# 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 nonexecuting processes saved in memory

# Multiprocessing: The (Traditional) Reality



■ **Save current registers in memory**

# Multiprocessing: The (Traditional) Reality

**Memory**

| Stack |
| Heap |
| Data |
| Code |

| Saved registers |

| Stack |
| Heap |
| Data |
| Code |

**. . .**

| Stack |
| Heap |
| Data |
| Code |

| Saved registers |

**CPU**

| Registers |

- Schedule next process for execution

# Multiprocessing: The (Traditional) Reality

**Memory**

| Stack | | Stack | | Stack |
|-------|---|-------|---|-------|
| Heap | | Heap | | Heap |
| Data | ... | Data | | Data |
| Code | | Code | | Code |
| Saved registers | | Saved registers | | Saved registers |

**CPU**

Registers

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

# Multiprocessing: The (Traditional) Reality

**Memory**

| Stack | Stack | Stack |
|-------|-------|-------|
| Heap | Heap | Heap |
| Data | Data | Data |
| Code | Code | Code |
| | **...** | Saved registers |

**CPU** — Registers

**CPU** — Registers

■ **Multicore processors**

▪ Multiple CPUs on single chip

▪ Share main memory (and some 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 running in parallel with each other**

Process A      Process B      Process C

**Time**

# 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***

# Execrise

- **P509 Textbook (Chinese Version)**
- **Execrise 8.1**

# 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(-1);
}
```

# 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(-1);
}
```

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

- **But, must think about application.  Not always appropriate to exit when something goes wrong.**

# 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:**
  - Receiving a signal whose default action is to terminate (next lecture)
  - Returning from the `main` routine
  - Calling the `exit` function

- `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***

# Conceptual View of `fork`



- **Make complete copy of execution state**
  - Designate one as parent and one as child
  - Resume execution of parent or child

# fork Example

```c
int main(int argc, char** argv)
{
    pid_t pid;
    int x = 1;

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

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

- **Call once, return twice**
- **Concurrent execution**
  - Can't predict execution order of parent and child

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

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

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

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

# Execrise

- P516 Textbook (Chinese Version)
- Execrise 8.2

# Making `fork` More Nondeterministic

- **Problem**
  - Linux scheduler does not create much run-to-run variance
  - Hides potential race conditions in nondeterministic programs
    - E.g., does `fork` return to child first, or to parent?
- **Solution**
  - Create custom version of library routine that inserts random delays along different branches
    - E.g., for parent and child in `fork`
  - Use runtime interpositioning to have program use special version of library code

# Variable delay `fork`

```c
/* fork wrapper function */
pid_t fork(void) {
    initialize();
    int parent_delay = choose_delay();
    int child_delay = choose_delay();
    pid_t parent_pid = getpid();
    pid_t child_pid_or_zero = real_fork();
    if (child_pid_or_zero > 0) {
        /* Parent */
        if (verbose) {
            printf("Fork.  Child pid=%d, delay = %dms.
                    Parent pid=%d, delay = %dms\n",
                    child_pid_or_zero, child_delay,
                    parent_pid, parent_delay);
            fflush(stdout);
        }
        ms_sleep(parent_delay);
    } else {
        /* Child */
        ms_sleep(child_delay);
    }
    return child_pid_or_zero;
}
```

*myfork.c*

# **`forkx2` Example**

```c
int main(int argc, char** argv)
{

    pid_t pid;
    int x = 1;

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

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

```
linux> ./fork2
parent: x=0
parent: x=-1
child : x=2
child : x=3
```

- **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 -> 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
  - two events in the graph has a fixed happening order

# Process Graph Example

```
int main(int argc, char** argv)
{
    pid_t pid;
    int x = 1;

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

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

# Interpreting Process Graphs

- **Original graph:**



```
                child: x=2
                    •────────────────────►•
                  printf                  exit

         x==1        parent: x=0
    •────────►•────────►•────────────────►•
   main      for      printf              exit
              k
```

- **Relabled graph:**



```
            ┌──────────►•────────────────►•
            │           e                  f
    •───────┴──►•────────►•────────────────►•
    a          b          c                  d
```

**Feasible total ordering:**



a   b   e   c   f   d

**Infeasible total ordering:**



a   b   f   c   e   d

# fork Example: Two consecutive forks

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



| Feasible | Infeasible |
|----------|------------|
| L0 | L0 |
| L1 | Bye |
| Bye | L1 |
| Bye | Bye |
| L1 | L1 |
| Bye | Bye |
| 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 | Infeasible |
|----------|------------|
| L0 | L0 |
| L1 | Bye |
| Bye | L1 |
| Bye | Bye |
| L2 | 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 | Infeasible |
|----------|------------|
| L0 | L0 |
| Bye | Bye |
| L1 | L1 |
| L2 | Bye |
| Bye | 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

```c
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 ttyp9    00:00:00 tcsh
 6639 ttyp9    00:00:03 forks
 6640 ttyp9    00:00:00 forks <defunct>
 6641 ttyp9    00:00:00 ps
linux> kill 6639
[1]    Terminated
linux> ps
  PID TTY          TIME CMD
 6585 ttyp9    00:00:00 tcsh
 6642 ttyp9    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

```c
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 ttyp9    00:00:00 tcsh
 6676 ttyp9    00:00:06 forks
 6677 ttyp9    00:00:00 ps
linux> kill 6676
linux> ps
  PID TTY          TIME CMD
 6585 ttyp9    00:00:00 tcsh
 6678 ttyp9    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
  - Implemented as syscall

*Parent Process*          *Kernel code*

syscall
...
*Exception*

*Returns*

And, potentially other user processes, including a child of parent

# `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

# `wait`: Synchronizing with Children

```c
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(s):**

| HC  | HP  |
| --- | --- |
| HP  | HC  |
| CT  | CT  |
| Bye | 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

```c
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)

```c
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*

# Execrise

- ■ P518 Textbook, Exercise 8.3
- ■ P520 Textbook, Exercise 8.4

# `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

Bottom of stack

| |
|---|
| Null-terminated **environment variable** strings |
| Null-terminated **command-line arg** strings |
| |
| `envp[n]  == NULL` |
| `envp[n-1]` |
| … |
| `envp[0]` |
| `argv[argc] = NULL` |
| `argv[argc-1]` |
| … |
| `argv[0]` |
| |
| Stack frame for `libc_start_main` |
| Future stack frame for `main` |

`environ` (global var)

`envp` (in `%rdx`)

`argv` (in `%rsi`)

`argc` (in `%rdi`)

Top of stack

# execve Example

- **Execute** `"/bin/ls -lt /usr/include"` **in child process using current environment:**

```
                    ┌──────────────────────┐
                    │ envp[n] = NULL        │
                    ├──────────────────────┤
                    │ envp[n-1]            │──────────▶  "PWD=/usr/droh"
                    ├──────────────────────┤
                    │ …                     │
                    ├──────────────────────┤
environ ──────────▶ │ envp[0]              │──────────▶  "USER=droh"
                    └──────────────────────┘

                    ┌──────────────────────┐
                    │ myargv[argc] = NULL   │
                    ├──────────────────────┤
(argc == 3)         │ myargv[2]            │──────────▶  "/usr/include"
                    ├──────────────────────┤
                    │ myargv[1]            │──────────▶  "-lt"
                    ├──────────────────────┤
myargv ───────────▶ │ myargv[0]            │──────────▶  "/bin/ls"
                    └──────────────────────┘
```

```
  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 any single core
- 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

# Exception Handling in MIT JOS Lab
# Based on Intel 80386

```
 4    // Trap numbers
 5    // These are processor defined:
 6    #define T_DIVIDE      0          // divide error
 7    #define T_DEBUG       1          // debug exception
 8    #define T_NMI         2          // non-maskable interrupt
 9    #define T_BRKPT       3          // breakpoint
10    #define T_OFLOW       4          // overflow
11    #define T_BOUND       5          // bounds check
12    #define T_ILLOP       6          // illegal opcode
13    #define T_DEVICE      7          // device not available
14    #define T_DBLFLT      8          // double fault
15    /* #define T_COPROC   9 */       // reserved (not generated by recent processors)
16    #define T_TSS        10          // invalid task switch segment
17    #define T_SEGNP      11          // segment not present
18    #define T_STACK      12          // stack exception
19    #define T_GPFLT      13          // general protection fault
20    #define T_PGFLT      14          // page fault
21    /* #define T_RES     15 */       // reserved
22    #define T_FPERR      16          // floating point error
23    #define T_ALIGN      17          // aligment check
24    #define T_MCHK       18          // machine check
25    #define T_SIMDERR    19          // SIMD floating point error
26
27    // These are arbitrarily chosen, but with care not to overlap
28    // processor defined exceptions or interrupt vectors.
29    #define T_SYSCALL    48          // system call
30    #define T_DEFAULT    500         // catchall
31
32    // Hardware IRQ numbers. We receive these as (IRQ_OFFSET+IRQ_WHATEVER)
33    // IRQ_OFFSET is defined in kern/picirq.h = 32
34    #define IRQ_TIMER      0
35    #define IRQ_KBD        1
36    #define IRQ_SPURIOUS   7
37    #define IRQ_IDE       14
38    #define IRQ_ERROR     19
```

```
        IDT                    trapentry.S          trap.c

+---------------+
|   &handler1   |----------> handler1:        trap (struct Trapframe *tf)
|               |              // do stuff     {
|               |              call trap          // handle the exception/interrupt
|               |              // ...           }
+---------------+
|   &handler2   |----------> handler2:
|               |              // do stuff
|               |              call trap
|               |              // ...
+---------------+
        .
        .
        .
+---------------+
|   &handlerX   |----------> handlerX:
|               |              // do stuff
|               |              call trap
|               |              // ...
+---------------+
```

```c
14    /* The TRAPHANDLER macro defines a globally-visible function for handling
15     * a trap.  It pushes a trap number onto the stack, then jumps to _alltraps.
16     * Use TRAPHANDLER for traps where the CPU automatically pushes an error code.
17     */
18    #define TRAPHANDLER(name, num)                                          \
19            .globl name;            /* define global symbol for 'name' */  \
20            .type name, @function;  /* symbol type is function */          \
21            .align 2;               /* align function definition */        \
22            name:                   /* function starts here */             \
23            pushl $(num);                                                   \
24            jmp _alltraps
25
44    _alltraps:
45            pushl %ds
46            pushl %es
47            pushal
48
49            movw $GD_KD, %ax
50            movw %ax, %ds
51            movw %ax, %es
52
53            pushl %esp /* esp as an argument to trap*/
54            call trap
55            addl $4, %esp /* jump esp */
56
57            popal
58            popl %es
59            popl %ds
60            addl $8, %esp /* jump trap number and error c
61            iret
```

```
325    .data
326    .global vectors
327    vectors:
328            .long vector0
329            .long vector1
330            .long vector2
331            .long vector3
332            .long vector4
333            .long vector5

67     TRAPHANDLER_NOEC(vector0,    0)
68     TRAPHANDLER_NOEC(vector1,    1)
69     TRAPHANDLER_NOEC(vector2,    2)
70     TRAPHANDLER_NOEC(vector3,    3)
71     TRAPHANDLER_NOEC(vector4,    4)
72     TRAPHANDLER_NOEC(vector5,    5)
73     TRAPHANDLER_NOEC(vector6,    6)
74     TRAPHANDLER_NOEC(vector7,    7)
75     TRAPHANDLER      (vector8,    8)
76     TRAPHANDLER_NOEC(vector9,    9)
77     TRAPHANDLER      (vector10,  10)
```

```c
64   void
65   idt_init(void)
66   {
67           extern struct Segdesc gdt[];
68
69           // LAB 3: Your code here.
70           int i;
71
72           for (i = 0; i < 256; i++)
73                   SETGATE(idt[i], 0, GD_KT, vectors[i], 0)
74
75           // Difference between interrupt gate and trap gate:
76           // After transfering the control, an interrupt gate clears IF to disable interrupts
77           // however, trap gate does not change the IF flag
78           // There will be errors in handling interrput when what's currently running is an exception
79           // JOS only allows interrupt happens in user space, and forbid it in kernel space, when a
80           // timer interrupt arrives, it will run sched_yield and eax is currently the syscall number
81           // which is >0, so user panic at lib/syscall.c
82           SETGATE(idt[T_SYSCALL], 0, GD_KT, vectors[T_SYSCALL], 3)
83           SETGATE(idt[T_BRKPT], 0, GD_KT, vectors[T_BRKPT], 3)
84           SETGATE(idt[T_OFLOW], 0, GD_KT, vectors[T_OFLOW], 3);
85           SETGATE(idt[T_BOUND], 0, GD_KT, vectors[T_BOUND], 3);
86
87           // Setup a TSS so that we get the right stack
88           // when we trap to the kernel.
89           ts.ts_esp0 = KSTACKTOP;
90           ts.ts_ss0 = GD_KD;
91
92           // Initialize the TSS field of the gdt.
93           gdt[GD_TSS >> 3] = SEG16(STS_T32A, (uint32_t) (&ts),
94                                   sizeof(struct Taskstate), 0);
95           gdt[GD_TSS >> 3].sd_s = 0;
96
97           // Load the TSS
98           ltr(GD_TSS);
99
100          // Load the IDT
101          asm volatile("lidt idt_pd");
```

```
185  void
186  trap(struct Trapframe *tf)
187  {
188          if ((tf->tf_cs & 3) == 3) {
189                  // Trapped from user mode.
190                  // Copy trap frame (which is currently on the stack)
191                  // into 'curenv->env_tf', so that running the environment
192                  // will restart at the trap point.
193                  assert(curenv);
194                  curenv->env_tf = *tf;
195                  // The trapframe on the stack should be ignored from here on.
196                  tf = &curenv->env_tf;
197          }
198
199          // Dispatch based on what type of trap occurred
200          trap_dispatch(tf);
201
202          // If we made it to this point, then no other environment was
203          // scheduled, so we should return to the current environment
204          // if doing so makes sense.
205          if (curenv && curenv->env_status == ENV_RUNNABLE)
206                  env_run(curenv);
207          else
208                  sched_yield();
209  }
```

```c
static void
trap_dispatch(struct Trapframe *tf)
{
        // Handle processor exceptions.
        // LAB 3: Your code here.
        struct PushRegs *regs;
        switch (tf->tf_trapno) {
        case T_SYSCALL:
                regs = &(tf->tf_regs);
                regs->reg_eax = syscall(regs->reg_eax, regs->reg_edx,
                        regs->reg_ecx, regs->reg_ebx, regs->reg_edi, regs->reg_esi);
                return;
        case T_PGFLT:
                page_fault_handler(tf);
                return;
        case T_BRKPT:
                monitor(tf);
                return;
        case T_DEBUG:
                monitor(tf);
                return;

        // Handle clock interrupts.
        // LAB 4: Your code here.
        case IRQ_OFFSET + IRQ_TIMER:
        // Add time tick increment to clock interrupts.
        // LAB 6: Your code here.
                time_tick();
                sched_yield();
                return;
        }
```

```c
471    // Dispatches to the correct kernel function, passing the arguments.
472    int32_t
473    syscall(uint32_t syscallno, uint32_t a1, uint32_t a2, uint32_t a3, uint32_t a4, uint32_t a5)
474    {
475            // Call the function corresponding to the 'syscallno' parameter.
476            // Return any appropriate return value.
477            // LAB 3: Your code here.
478
479            int32_t ret = 0;
480
481            switch (syscallno) {
482            case SYS_cputs:
483                    sys_cputs((char *)a1, (size_t)a2);
484                    break;
485            case SYS_cgetc:
486                    ret = sys_cgetc();
487                    break;
488            case SYS_getenvid:
489                    ret = sys_getenvid();
490                    break;
491            case SYS_env_destroy:
492                    ret = sys_env_destroy((envid_t)a1);
493                    break;
494            case SYS_yield:
495                    sys_yield();
496                    break;
497            case SYS_exofork:
498                    ret = sys_exofork();
499                    break;
500            case SYS_env_set_status:
501                    ret = sys_env_set_status((envid_t)a1, a2);
502                    break;
```