

POSIX process API

last time

kernel part of context switch

- save all registers; restore all registers

- trick: function calls save some registers automatically

user registers: save/restore on mode switch

- part of exception handling (even if no context switch)

thread + process control blocks

[3:30pm] myproc() as processor-local “variable”

process control block

some data structure needed to represent a process

called **Process Control Block**

process control block

some data structure needed to represent a process

called **Process Control Block**

xv6: `struct proc`

xv6: struct proc

```
struct proc {  
    uint sz;                // Size of process memory (bytes)  
    pde_t* pgdir;           // Page table  
    char *kstack;           // Bottom of kernel stack for this process  
    enum procstate state;   // Process state  
    int pid;                // Process ID  
    struct proc *parent;    // Parent process  
    struct trapframe *tf;   // Trap frame for current syscall  
    struct context *context; // swtch() here to run process  
    void *chan;             // If non-zero, sleeping on chan  
    int killed;             // If non-zero, have been killed  
    struct file *ofile[NOFILE]; // Open files  
    struct inode *cwd;      // Current directory  
    char name[16];         // Process name (debugging)  
};
```

xv6: struct proc

pointers to current registers/PC of process (user and kernel)
stored on its kernel stack
(if not currently running)

```
struct proc {  
    uint sz;  
    pde_t* pg;  
    char *kstack;  
    enum proc_state state; // thread's state  
    int pid;                // Process ID  
    struct proc *parent;    // Parent process  
    struct trapframe *tf;   // Trap frame for current syscall  
    struct context *context; // swtch() here to run process  
    void *chan;             // If non-zero, sleeping on chan  
    int killed;              // If non-zero, have been killed  
    struct file *ofile[NOFILE]; // Open files  
    struct inode *cwd;       // Current directory  
    char name[16];          // Process name (debugging)  
};
```

SS

xv6: struct proc

```
struct proc {  
    uint sz;  
    pde_t* pgdir;  
    char *kstack;  
    enum procstate state;  
    int pid;  
    struct proc *parent;  
    struct trapframe *tf;  
    struct context *context;  
    void *chan;  
    int killed;  
    struct file *ofile[NOFILE];  
    struct inode *cwd;  
    char name[16];  
};
```

the kernel stack for this process
every process has one kernel stack

```
// Size of process memory (bytes)  
// Page table  
// Bottom of kernel stack for this process  
// Process state  
// Process ID  
// Parent process  
// Trap frame for current syscall  
// swtch() here to run process  
// If non-zero, sleeping on chan  
// If non-zero, have been killed  
// Open files  
// Current directory  
// Process name (debugging)
```

xv6: struct proc

```
struct proc {  
    uint sz;  
    pde_t* pgtable;  
    char *kstack;  
    enum procstate state;  
    int pid;  
    struct proc *parent;  
    struct trapframe *tf;  
    struct context *context;  
    void *chan;  
    int killed;  
    struct file *ofile[NOFILE];  
    struct inode *cwd;  
    char name[16];  
};  
  
enum procstate {  
    UNUSED, EMBRYO, SLEEPING,  
    RUNNABLE, RUNNING, ZOMBIE  
};
```

is process running?
or waiting?
or finished?
if waiting,
waiting for what (chan)?

*// Process state
// Process ID
// Parent process
// Trap frame for current syscall
// swtch() here to run process
// If non-zero, sleeping on chan
// If non-zero, have been killed
// Open files
// Current directory
// Process name (debugging)*

xv6: struct proc

```
struct proc {  
    uint sz;  
    pde_t* pgdir;  
    char *kstack;  
    enum procstate state;  
    int pid;  
    struct proc *parent;  
    struct trapframe *tf;  
    struct context *context;  
    void *chan;  
    int killed;  
    struct file *ofile[NOFILE];  
    struct inode *cwd;  
    char name[16];  
};
```

process ID

to identify process in system calls

```
// Size of process memory (bytes)  
// Page table  
// Bottom of kernel stack for this process  
// Process state  
// Process ID  
// Parent process  
// Trap frame for current syscall  
// swtch() here to run process  
// If non-zero, sleeping on chan  
// If non-zero, have been killed  
// Open files  
// Current directory  
// Process name (debugging)
```

xv6: struct proc

```
struct proc {  
    uint sz;  
    pde_t* pgdir;  
    char *kstack;  
    enum procstate state;  
    int pid;  
    struct proc *parent;  
    struct trapframe *tf;  
    struct context *context;  
    void *chan;  
    int killed;  
    struct file *ofile[NOFILE];  
    struct inode *cwd;  
    char name[16];  
};
```

// Size of process memory (bytes)
// Page table
// Bottom of kernel stack for this process
// Process state
// Proc
// Pare
// Trap
// swtc
// If n
// If non-zero, have been killed
// Open files
// Current directory
// Process name (debugging)

information about address space
pgdir — used by processor
sz — used by OS only

xv6: struct proc

information about open files, etc.

```
struct proc {  
    uint sz;  
    pde_t* pgdir;  
    char *kstack;  
    enum procstate state;  
    int pid;  
    struct proc *parent;  
    struct trapframe *tf;  
    struct context *context;  
    void *chan;  
    int killed;  
    struct file *ofile[NOFILE];  
    struct inode *cwd;  
    char name[16];  
};
```

// Size of process memory (bytes)
// Page table
// Bottom of kernel stack for this process
// Process state
// Process ID
// Parent process
// Trap frame for current syscall
// swtch() here to run process
// If non-zero, sleeping on chan
// If non-zero, have been killed
// Open files
// Current directory
// Process name (debugging)

process control blocks generally

contains process's context(s) (registers, PC, ...)

- if context is not on a CPU

- (in xv6: pointers to these, actual location: process's kernel stack)

process's status — running, waiting, etc.

information for system calls, etc.

- open files

- memory allocations

- process IDs

- related processes

xv6 myproc

xv6 function: `myproc()`

retrieves pointer to currently running struct `proc`

myproc: using a global variable

```
struct cpu cpus[NCPU];
```

```
struct proc*  
myproc(void) {  
    struct cpu *c;  
    ...  
    c = mycpu();    /* finds entry of cpus array  
                     using special "ID" register  
                     as array index */  
  
    p = c->proc;  
    ...  
    return p;  
}
```

this class: focus on Unix

Unix-like OSes will be our focus

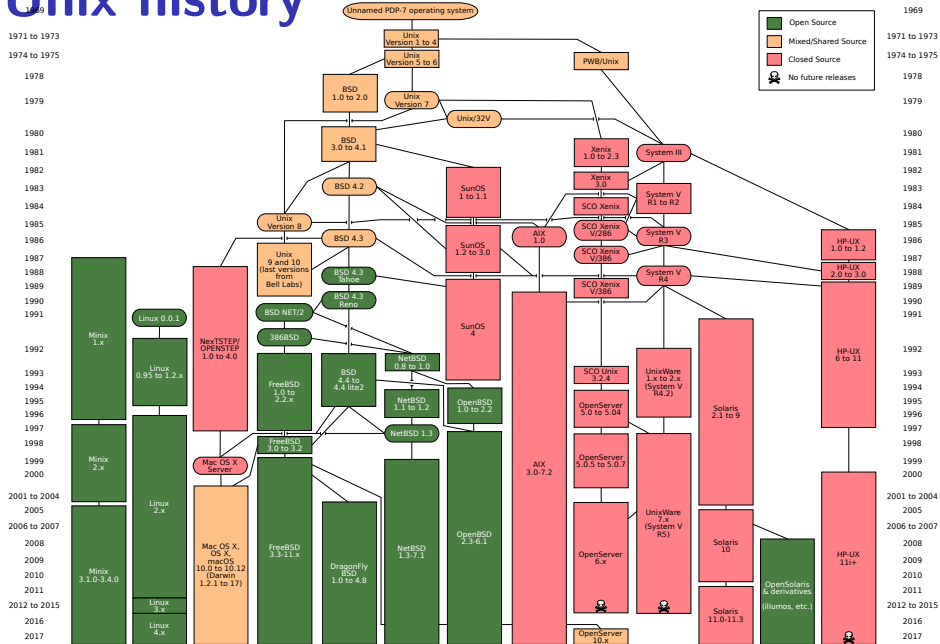
we have source code

used to from 2150, etc.?

have been around for a while

xv6 imitates Unix

Unix history



POSIX: standardized Unix

Portable Operating System Interface (POSIX)

“standard for Unix”

current version online:

<https://pubs.opengroup.org/onlinepubs/9699919799/>

(almost) followed by most current Unix-like OSes

...but OSes add extra features

...and POSIX doesn't specify everything

what POSIX defines

POSIX specifies the **library and shell interface**
source code compatibility

doesn't care what is/is not a system call...

doesn't specify binary formats...

idea: write applications for POSIX, recompile and run on all implementations

this was a very important goal in the 80s/90s
at the time, no dominant Unix-like OS (Linux was very immature)

POSIX process management

essential operations

process information: `getpid`

process creation: `fork`

running programs: `exec*`

also `posix_spawn` (not widely supported), ...

waiting for processes to finish: `waitpid` (or `wait`)

process destruction, 'signaling': `exit`, `kill`

POSIX process management

essential operations

process information: `getpid`

process creation: `fork`

running programs: `exec*`

also `posix_spawn` (not widely supported), ...

waiting for processes to finish: `waitpid` (or `wait`)

process destruction, 'signaling': `exit`, `kill`

getpid

```
pid_t my_pid = getpid();  
printf("my pid is %ld\n", (long) my_pid);
```

process ids in ps

```
cr4bd@machine:~$ ps
```

PID	TTY	TIME	CMD
14777	pts/3	00:00:00	bash
14798	pts/3	00:00:00	ps

POSIX process management

essential operations

process information: `getpid`

process creation: `fork`

running programs: `exec*`

also `posix_spawn` (not widely supported), ...

waiting for processes to finish: `waitpid` (or `wait`)

process destruction, 'signaling': `exit`, `kill`

fork

`pid_t fork()` — copy the current process

returns twice:

in *parent* (original process): pid of new *child* process

in *child* (new process): 0

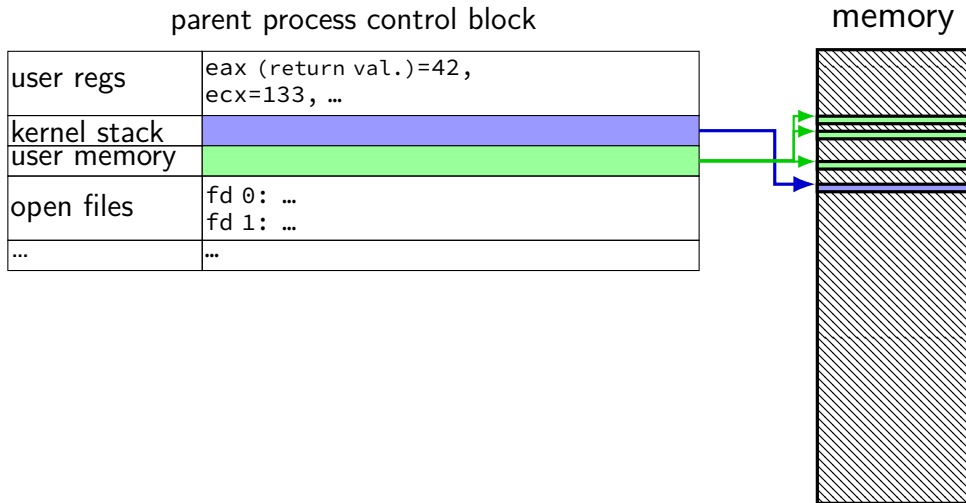
everything (but pid) duplicated in parent, child:

memory

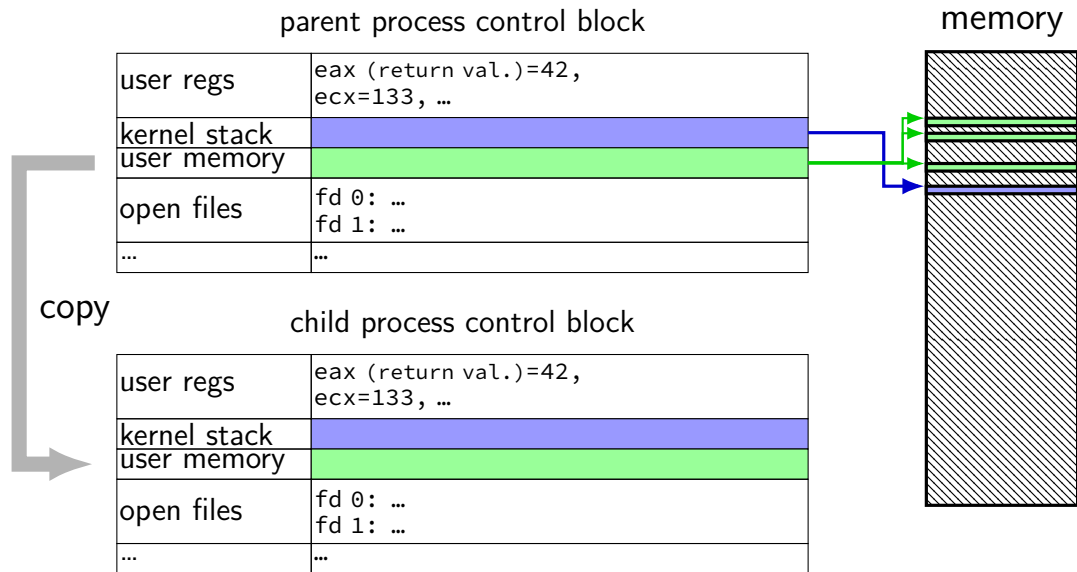
file descriptors (later)

registers

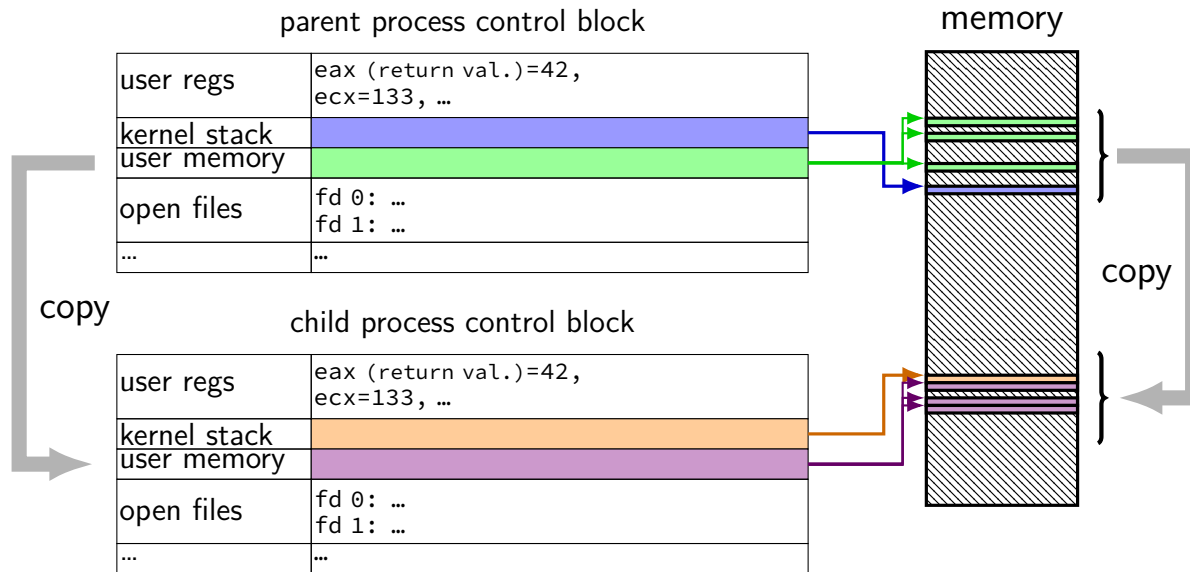
fork and PCBs



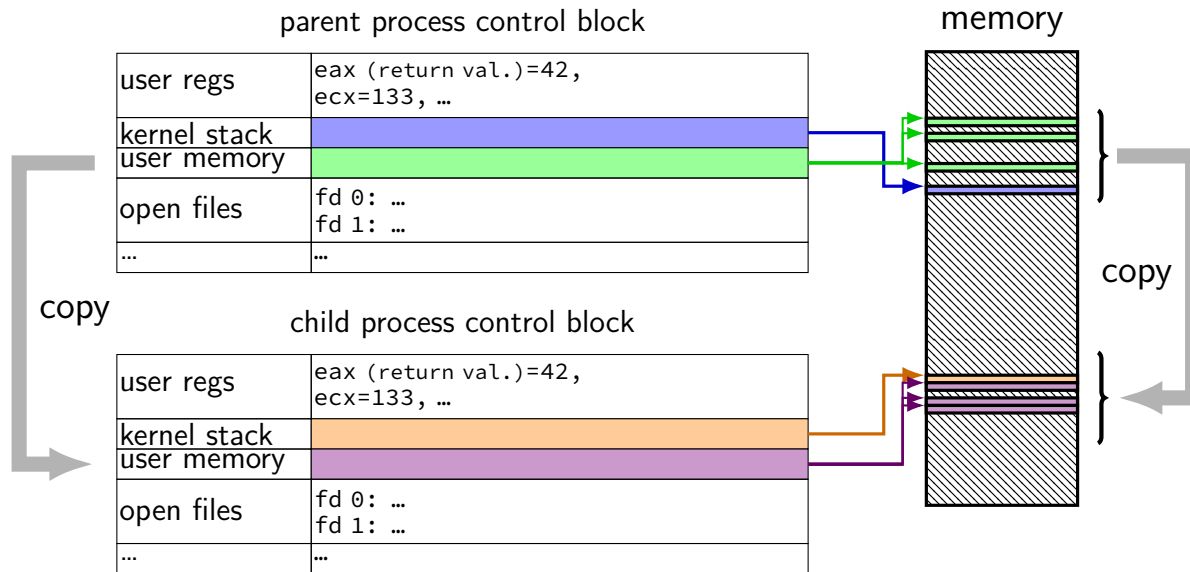
fork and PCBs



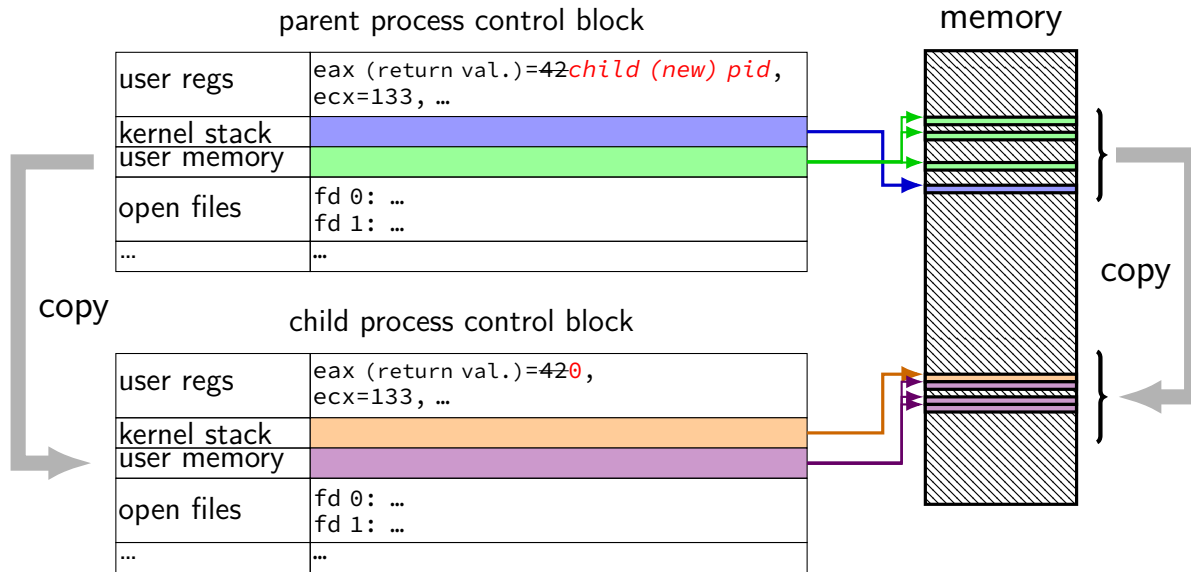
fork and PCBs



fork and PCBs



fork and PCBs



fork example

```
#include <stdlib.h>
#include <stdio.h>
#include <unistd.h>
#include <sys/types.h>
int main(int argc, char *argv[]) {
    pid_t pid = getpid();
    printf("Parent pid: %d\n", (int) pid);
    pid_t child_pid = fork();
    if (child_pid > 0) {
        /* Parent Process */
        pid_t my_pid = getpid();
        printf("[%d] parent of [%d]\n", (int) my_pid, (int) child_pid);
    } else if (child_pid == 0) {
        /* Child Process */
        pid_t my_pid = getpid();
        printf("[%d] child\n", (int) my_pid);
    } else {
        perror("Fork failed");
    }
    return 0;
}
```

fork example

```
#include <stdlib.h>
#include <stdio.h>
#include <unistd.h>
#include <sys/types.h>
int main(int argc, char *argv[]) {
    pid_t pid = getpid();
    printf("Parent pid: %d\n", (int) pid);
    pid_t child_pid = fork();
    if (child_pid > 0) {
        /* Parent Process */
        pid_t my_pid = getpid();
        printf("[%d] parent of [%d]\n", (int) my_pid, (int) child_pid);
    } else if (child_pid == 0) {
        /* Child Process */
        pid_t my_pid = getpid();
        printf("[%d] child\n", (int) my_pid);
    } else {
        perror("Fork failed");
    }
    return 0;
}
```

getpid — returns current process pid

fork example

```
#include <stdlib.h>
#include <stdio.h>
#include <unistd.h>
#include <sys/types.h>
int main(int argc, char **argv) {
    pid_t pid = fork();
    printf("Parent PID: %d\n", pid);
    pid_t child_pid = fork();
    if (child_pid > 0) {
        /* Parent Process */
        pid_t my_pid = getpid();
        printf("[%d] parent of [%d]\n", (int) my_pid, (int) child_pid);
    } else if (child_pid == 0) {
        /* Child Process */
        pid_t my_pid = getpid();
        printf("[%d] child\n", (int) my_pid);
    } else {
        perror("Fork failed");
    }
    return 0;
}
```

cast in case pid_t isn't int
POSIX doesn't specify (some systems it is, some not...)
(not necessary if you were using C++'s cout, etc.)

fork example

```
#include <stdlib.h>
#include <stdio.h>
#include <unistd.h>
#include <errno.h>
int main()
{
    pid_t child_pid;
    printf("Forking...\n");
    child_pid = fork();
    if (child_pid > 0) {
        /* Parent Process */
        pid_t my_pid = getpid();
        printf("[%d] parent of [%d]\n", (int) my_pid, (int) child_pid);
    } else if (child_pid == 0) {
        /* Child Process */
        pid_t my_pid = getpid();
        printf("[%d] child\n", (int) my_pid);
    } else {
        perror("Fork failed");
    }
    return 0;
}
```

prints out Fork failed: *error message*
(example *error message*: "Resource temporarily unavailable")
from error number stored in special global variable `errno`

fork example

```
#include <stdlib.h>
#include <stdio.h>
#include <unistd.h>
#include <sys/types.h>
int main(int argc, char *argv[]) {
    pid_t pid = getpid();
    printf("Parent pid: %d\n", (int) pid);
    pid_t child_pid = fork();
    if (child_pid > 0) {
        /* Parent Process */
        pid_t my_pid = getpid();
        printf("[%d] parent of [%d]\n", (int) my_pid, (int) child_pid);
    } else if (child_pid == 0) {
        /* Child Process */
        pid_t my_pid = getpid();
        printf("[%d] child\n", (int) my_pid);
    } else {
        perror("Fork failed");
    }
    return 0;
}
```

Example output:

Parent pid: 100

[100] parent of [432]

[432] child

a fork question

```
int main() {  
    pid_t pid = fork();  
    if (pid == 0) {  
        printf("In child\n");  
    } else {  
        printf("Child %d\n", pid);  
    }  
    printf("Done!\n");  
}
```

Exercise: Suppose the pid of the parent process is 99 and child is 100. Give **two** possible outputs. (Assume no crashes, etc.)

a fork question

```
int main() {  
    pid_t pid = fork();  
    if (pid == 0) {  
        printf("In child\n");  
    } else {  
        printf("Child %d\n", pid);  
    }  
    printf("Done!\n");  
}
```

Exercise: Suppose the pid of the parent process is 99 and child is 100. Give **two** possible outputs. (Assume no crashes, etc.)



Child 100
In child
Done!
Done!



In child
Done!
Child 100
Done!

POSIX process management

essential operations

process information: `getpid`

process creation: `fork`

running programs: `exec*`

also `posix_spawn` (not widely supported), ...

waiting for processes to finish: `waitpid` (or `wait`)

process destruction, 'signaling': `exit`, `kill`

exec*

exec* — **replace** current program with new program

* — multiple variants

same pid, new process image

```
int execl(const char *path, const char **argv)
```

path: new program to run

argv: array of arguments, terminated by null pointer

execv example

```
...
child_pid = fork();
if (child_pid == 0) {
    /* child process */
    char *args[] = {"ls", "-l", NULL};
    execv("/bin/ls", args);
    /* execv doesn't return when it works.
    So, if we got here, it failed. */
    perror("execv");
    exit(1);
} else if (child_pid > 0) {
    /* parent process */
    ...
}
```

execv example

```
...
child_pid = fork();
if (child_pid == 0) {
    /* child process */
    char *args[] = {"ls", "-l", NULL};
    execv("/bin/ls", args);
    /* execv doesn't return when it works.
        So, if we got
    perror("execv");
    exit(1);
} else if (child_pid > 0) {
    /* parent process */
    ...
}
```

used to compute argv, argc
when program's main is run

convention: first argument is program name

execv example

```
...
child_pid = fork();
if (child_pid == 0) {
    /* child process */
    char *args[] = {"ls", "-l", NULL};
    execv("/bin/ls", args);
    /* execv doesn't return when it works.
```

```
    So, if we got here,
    perror("execv");
    exit(1);
} else if (child_pid > 0) {
    /* parent process */
```

```
    ...
}
```

path of executable to run
need not match first argument
(but probably should match it)

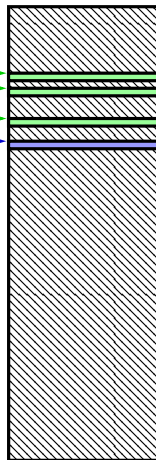
on Unix /bin is a directory
containing many common programs,
including ls ('list directory')

exec and PCBs

the process control block

user regs	eax=42, ecx=133, ...
kernel stack	
user memory	
open files	fd 0: (terminal ...) fd 1: ...
...	...

memory

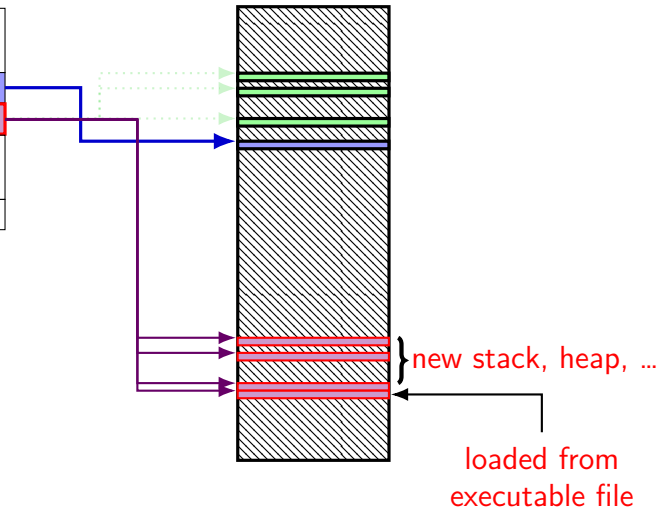


exec and PCBs

the process control block

user regs	eax=42 init. val. , ecx=133 init. val. , ...
kernel stack	
user memory	
open files	fd 0: (terminal ...) fd 1: ...
...	...

memory

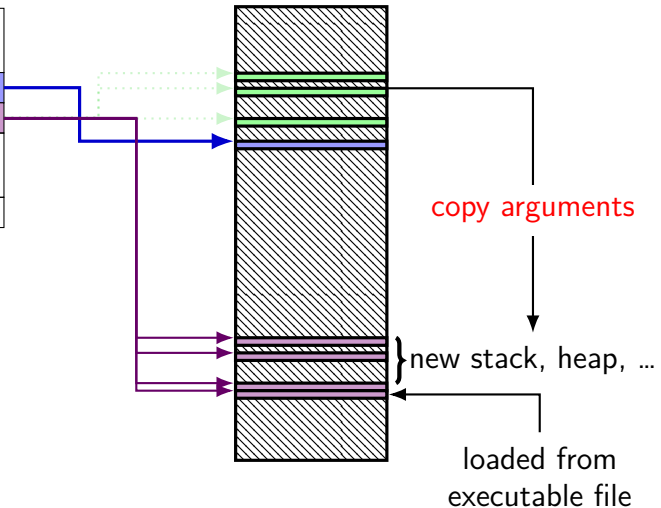


exec and PCBs

the process control block

user regs	<code>eax=42</code> <i>init. val.</i> , <code>ecx=133</code> <i>init. val.</i> , ...
kernel stack	
user memory	
open files	<code>fd 0:</code> (terminal ...) <code>fd 1:</code> ...
...	...

memory



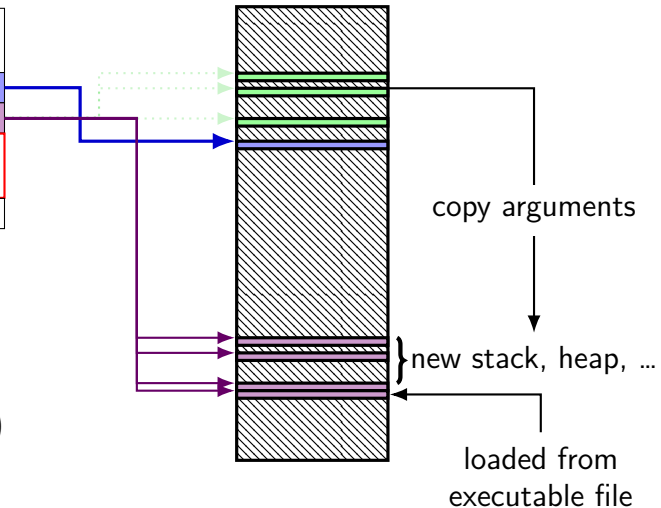
exec and PCBs

the process control block

user regs	<code>eax=42</code> <i>init. val.</i> , <code>ecx=133</code> <i>init. val.</i> , ...
kernel stack	
user memory	
open files	<code>fd 0: (terminal ...)</code> <code>fd 1: ...</code>
...	...

not changed!
(more on this later)

memory



exec and PCBs

the process control block

user regs	<code>eax=42</code> <i>init. val.</i> , <code>ecx=133</code> <i>init. val.</i> , ...
kernel stack	
user memory	
open files	<code>fd 0:</code> (terminal ...) <code>fd 1:</code> ...
...	...

not changed!
(more on this later)

memory

old memory
discarded

copy arguments

} new stack, heap, ...

loaded from
executable file

why fork/exec?

could just have a function to spawn a new program

Windows `CreateProcess()`; POSIX's (rarely used) `posix_spawn`

some other OSs do this (e.g. Windows)

needs to include API to set new program's state

e.g. without fork: need function to set new program's current directory

e.g. with fork: just change your current directory before exec

but allows OS to avoid 'copy everything' code

probably makes OS implementation easier

posix_spawn

```
pid_t new_pid;  
const char argv[] = { "ls", "-l", NULL };  
int error_code = posix_spawn(  
    &new_pid,  
    "/bin/ls",  
    NULL /* null = copy current process's open files;  
          if not null, do something else */,  
    NULL /* null = no special settings for new process */,  
    argv,  
    NULL /* null = copy current process's "environment variables";  
          if not null, do something else */  
);  
if (error_code == 0) {  
    /* handle error */  
}
```


some opinions (via HotOS '19)

A fork() in the road

Andrew Baumann
Microsoft Research

Jonathan Appavoo
Boston University

Orran Krieger
Boston University

Timothy Roscoe
ETH Zurich

ABSTRACT

The received wisdom suggests that Unix's unusual combination of `fork()` and `exec()` for process creation was an inspired design. In this paper, we argue that `fork` was a clever hack for machines and programs of the 1970s that has long outlived its usefulness and is now a liability. We catalog the ways in which `fork` is a terrible abstraction for the modern programmer to use, describe how it compromises OS implementations, and propose alternatives.

POSIX process management

essential operations

process information: `getpid`

process creation: `fork`

running programs: `exec*`

also `posix_spawn` (not widely supported), ...

waiting for processes to finish: `waitpid` (or `wait`)

process destruction, 'signaling': `exit`, `kill`

wait/waitpid

```
pid_t waitpid(pid_t pid, int *status,  
              int options)
```

wait for a child process (with `pid=pid`) to finish

sets `*status` to its “status information”

`pid=-1` → wait for any child process instead

options? see manual page (command `man waitpid`)

0 — no options

exit statuses

```
int main() {  
    return 0;  /* or exit(0); */  
}
```

waitpid example

```
#include <sys/wait.h>
...
child_pid = fork();
if (child_pid > 0) {
    /* Parent process */
    int status;
    waitpid(child_pid, &status, 0);
} else if (child_pid == 0) {
    /* Child process */
    ...
}
```

the status

```
#include <sys/wait.h>

...
waitpid(child_pid, &status, 0);
if (WIFEXITED(status)) {
    printf("main returned or exit called with %d\n",
           WEXITSTATUS(status));
} else if (WIFSIGNALED(status)) {
    printf("killed by signal %d\n", WTERMSIG(status));
} else {
    ...
}
```

“status code” encodes both return value and if exit was abnormal
W* macros to decode it

the status

```
#include <sys/wait.h>

...
waitpid(child_pid, &status, 0);
if (WIFEXITED(status)) {
    printf("main returned or exit called with %d\n",
           WEXITSTATUS(status));
} else if (WIFSIGNALED(status)) {
    printf("killed by signal %d\n", WTERMSIG(status));
} else {
    ...
}
```

“status code” encodes both return value and if exit was abnormal
W* macros to decode it

aside: signals

signals are a way of communicating between processes

they are also how abnormal termination happens

kernel communicating “something bad happened” → kills program by default

wait's status will tell you when and what signal killed a program

constants in signal.h

SIGINT — control-C

SIGTERM — kill command (by default)

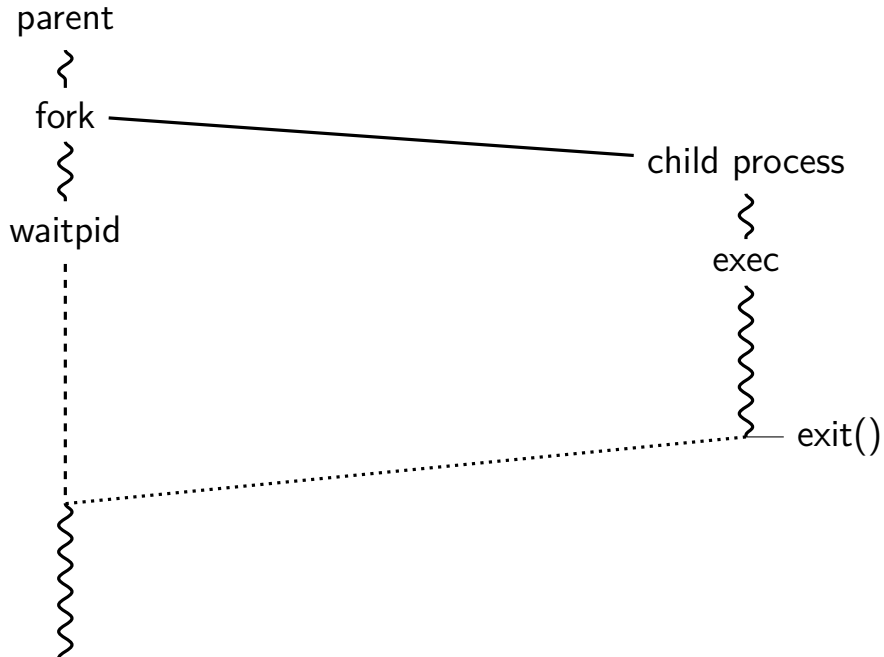
SIGSEGV — segmentation fault

SIGBUS — bus error

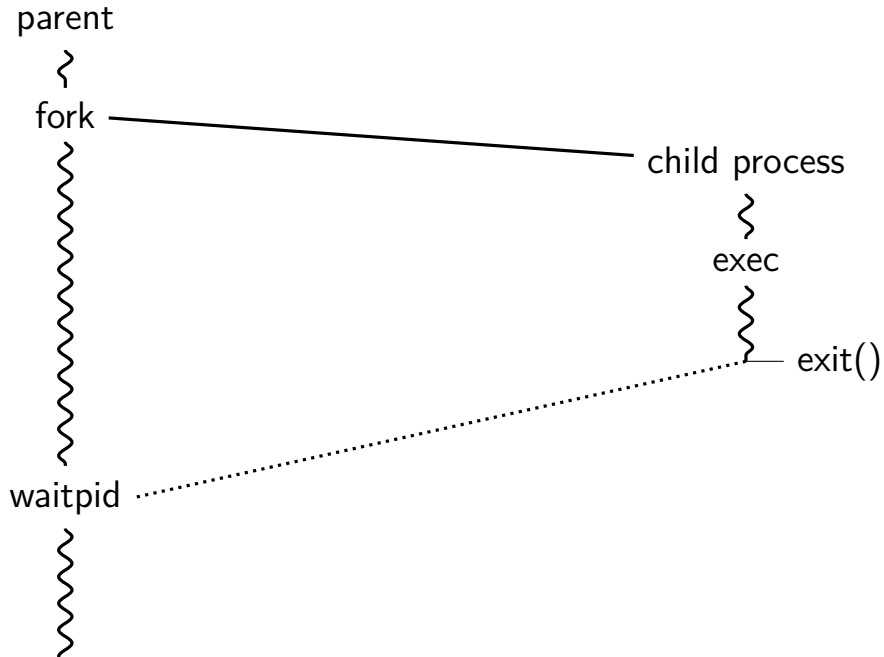
SIGABRT — abort() library function

...

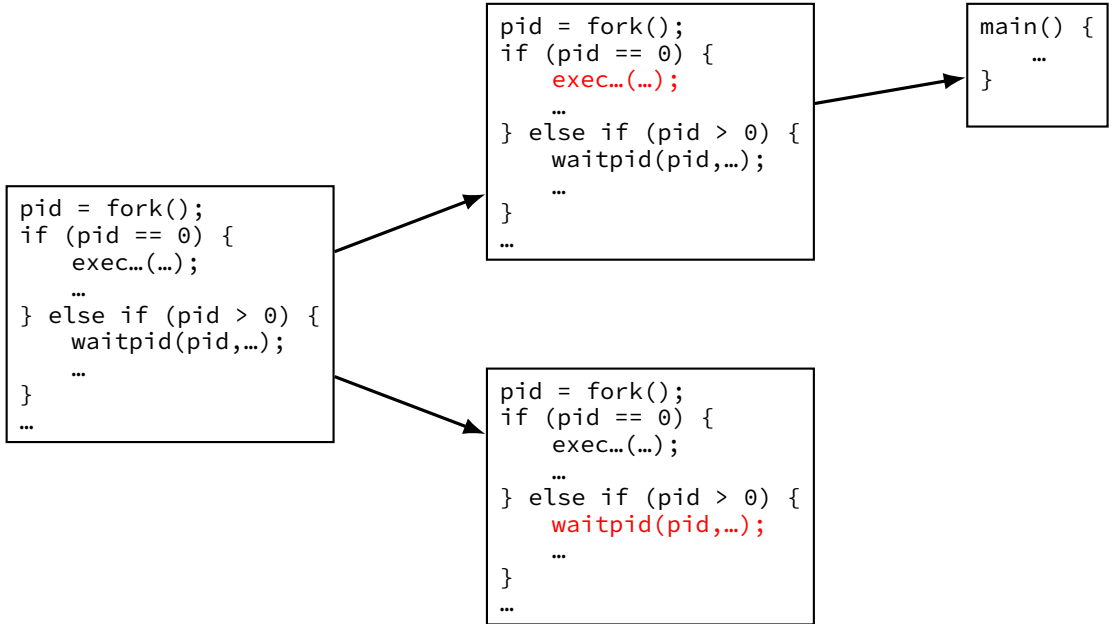
typical pattern



typical pattern (alt)



typical pattern (detail)



multiple processes?

```
while (...) {  
    pid = fork();  
    if (pid == 0) {  
        exec ...  
    } else if (pid > 0) {  
        pids.push_back(pid);  
    }  
}  
  
/* retrieve exit statuses in order */  
for (pid_t pid : pids) {  
    waitpid(pid, ...);  
    ...  
}
```

multiple processes?

```
while (...) {  
    pid = fork();  
    if (pid == 0) {  
        exec ...  
    } else if (pid > 0) {  
        pids.push_back(pid);  
    }  
}
```

```
/* retrieve exit statuses as processes finish */  
while ((pid = waitpid(-1, ...)) != -1) {  
    handleProcessFinishing(pid);  
}
```

POSIX process management

essential operations

process information: `getpid`

process creation: `fork`

running programs: `exec*`

also `posix_spawn` (not widely supported), ...

waiting for processes to finish: `waitpid` (or `wait`)

process destruction, 'signaling': `exit`, `kill`

backup slides

aside: environment variables (1)

key=value pairs associated with every process:

```
$ printenv
```

```
MODULE_VERSION_STACK=3.2.10
```

```
MANPATH=/opt/puppetlabs/puppet/share/man
```

```
XDG_SESSION_ID=754
```

```
HOSTNAME=labsrv01
```

```
SELINUX_ROLE_REQUESTED=
```

```
TERM=screen
```

```
SHELL=/bin/bash
```

```
HISTSIZE=1000
```

```
SSH_CLIENT=128.143.67.91 58432 22
```

```
SELINUX_USE_CURRENT_RANGE=
```

```
QTDIR=/usr/lib64/qt-3.3
```

```
OLDPWD=/zf14/cr4bd
```

```
QTINC=/usr/lib64/qt-3.3/include
```

```
SSH_TTY=/dev/pts/0
```

```
QT_GRAPHICSSYSTEM_CHECKED=1
```

```
USER=cr4bd
```

```
LS_COLORS=rs=0:di=01;34:ln=01;36:mh=00:pi=40;33:so=01;35:do=01;35:bd=40;33;01:cd=40;33;01:or=
```

```
MODULE_VERSION=3.2.10
```

```
MAIL=/var/spool/mail/cr4bd
```

```
PATH=/zf14/cr4bd/.cargo/bin:/zf14/cr4bd/bin:/usr/lib64/qt-3.3/bin:/usr/local/bin:/usr/bin:/us
```

```
PWD=/zf14/cr4bd
```

```
LANG=en_US.UTF-8
```

```
MODULEPATH=/sw/centos/Modules/modulefiles:/sw/linux-any/Modules/modulefiles
```

```
LOADEDMODULES=
```

```
KDEDIRS=/usr
```


aside: environment variables (2)

environment variable library functions:

`getenv("KEY")` \rightarrow *value*

`putenv("KEY=value")` (sets KEY to *value*)

`setenv("KEY", "value")` (sets KEY to *value*)

```
int execve(char *path, char **argv, char **envp)
```

```
char *envp[] = { "KEY1=value1", "KEY2=value2", NULL };
```

```
char *argv[] = { "somecommand", "some arg", NULL };
```

```
execve("/path/to/somecommand", argv, envp);
```

normal exec versions — keep same environment variables

aside: environment variables (3)

interpretation up to programs, but common ones...

`PATH=/bin:/usr/bin`

to run a program 'foo', look for an executable in `/bin/foo`, then `/usr/bin/foo`

`HOME=/zf14/cr4bd`

current user's home directory is `'/zf14/cr4bd'`

`TERM=screen-256color`

your output goes to a 'screen-256color'-style terminal

...

waiting for all children

```
#include <sys/wait.h>

...
while (true) {
    pid_t child_pid = waitpid(-1, &status, 0);
    if (child_pid == (pid_t) -1) {
        if (errno == ECHILD) {
            /* no child process to wait for */
            break;
        } else {
            /* some other error */
        }
    }
    /* handle child_pid exiting */
}
```

'waiting' without waiting

```
#include <sys/wait.h>
```

```
...
```

```
pid_t return_value = waitpid(child_pid, &status, WNOHANG);  
if (return_value == (pid_t) 0) {  
    /* child process not done yet */  
} else if (child_pid == (pid_t) -1) {  
    /* error */  
} else {  
    /* handle child_pid exiting */  
}
```

running in background

```
$ ./long_computation >tmp.txt &  
[1] 4049  
$ ...  
[1]+  Done                  ./long_computation > tmp.txt  
$ cat tmp.txt  
the result is ...
```

& — run a program in “background”

initially output PID (above: 4049)

print out after terminated

one way: use `waitpid` with option saying “don’t wait”

execv and const

```
int execv(const char *path, char *const *argv);
```

argv is a pointer to constant pointer to char

probably should be a pointer to constant pointer to *constant* char

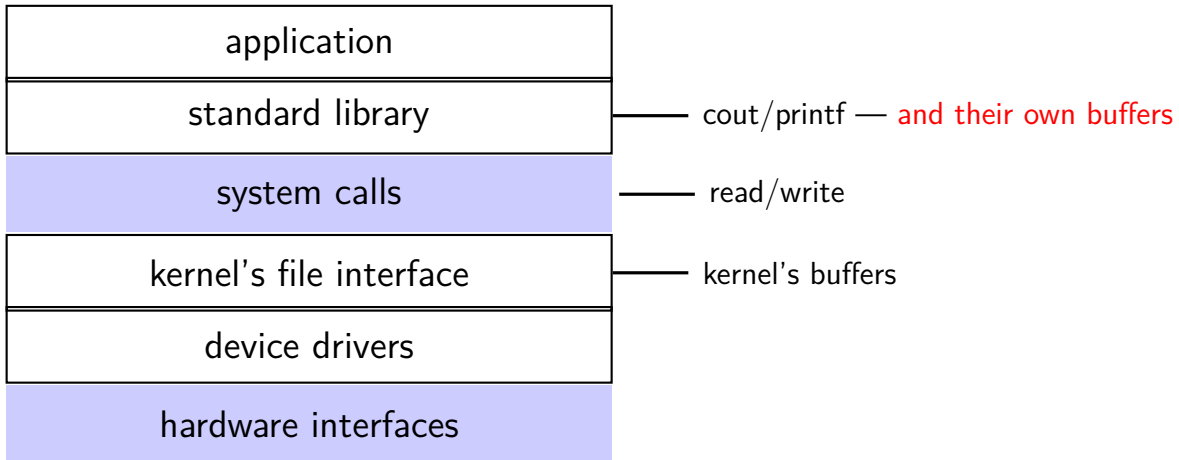
...this causes some awkwardness:

```
const char *array[] = { /* ... */ };  
execv(path, array); // ERROR
```

solution: cast

```
const char *array[] = { /* ... */ };  
execv(path, (char **) array); // or (char * const *)
```

layering



why the extra layer

better (but more complex to implement) interface:

- read line

- formatted input (scanf, cin into integer, etc.)

- formatted output

less system calls (bigger reads/writes) sometimes faster

- buffering can combine multiple in/out library calls into one system call

more portable interface

- cin, printf, etc. defined by C and C++ standards

parent and child processes

every process (but process id 1) has a *parent process* (getppid())

this is the process that can wait for it

creates tree of processes (Linux pstree command):

```
init(1)-+-ModemManager(919)-+-{ModemManager}(972)
|                               +-{ModemManager}(1064)
|                               |
|                               +-NetworkManager(1160)-+-dhclient(1755)
|                                                       |
|                                                       +-dnsmasq(1985)
|                                                       |
|                                                       +-{NetworkManager}(1180)
|                                                       |
|                                                       +-{NetworkManager}(1194)
|                                                       +-{NetworkManager}(1195)
|
| -accounts-daemon(1649)-+-{accounts-daemon}(1757)
|                       +-{accounts-daemon}(1758)
|
| -acpid(1338)
| -apache2(3165)-+-apache2(4125)-+-{apache2}(4126)
|                  +-{apache2}(4127)
|                  |
|                  +-apache2(28920)-+-{apache2}(28926)
|                  +-{apache2}(28960)
|                  |
|                  +-apache2(28921)-+-{apache2}(28927)
|                  +-{apache2}(28963)
|                  |
|                  +-apache2(28922)-+-{apache2}(28928)
|                  +-{apache2}(28961)
|                  |
|                  +-apache2(28923)-+-{apache2}(28930)
|                  +-{apache2}(28962)
|                  |
|                  +-apache2(28925)-+-{apache2}(28958)
|                  +-{apache2}(28965)
|                  |
|                  +-apache2(32165)-+-{apache2}(32166)
|                  +-{apache2}(32167)
|
| -at-spi-bus-laun(2252)-+-dbus-daemon(2269)
|                       |
|                       +-{at-spi-bus-laun}(2266)
|                       |
|                       +-{at-spi-bus-laun}(2268)
|                       |
|                       +-{at-spi-bus-laun}(2270)
|
| -at-spi2-registr(2275)-+-{at-spi2-registr}(2282)
|
| -atd(1633)
|
| -automount(13454)-+-{automount}(13455)
|                  +-{automount}(13456)
|                  |
|                  +-{automount}(13461)
|                  |
|                  +-{automount}(13464)
|                  +-{automount}(13465)
|
| -avahi-daemon(934)-+-avahi-daemon(944)
|
| -bluetoothd(924)
|
| -colord(1193)-+-{colord}(1329)
|               +-{colord}(1330)
|
| -mongod(1336)-+-{mongod}(1556)
|               +-{mongod}(1557)
|               |
|               +-{mongod}(1983)
|               |
|               +-{mongod}(2031)
|               |
|               +-{mongod}(2047)
|               |
|               +-{mongod}(2048)
|               |
|               +-{mongod}(2049)
|               |
|               +-{mongod}(2050)
|               |
|               +-{mongod}(2051)
|               +-{mongod}(2052)
|
| -mosh-server(19098)-+-bash(19091)---tmux(5442)
| -mosh-server(21996)---bash(21997)
| -mosh-server(22533)---bash(22534)---tmux(22588)
| -nn-applet(2580)-+-{nn-applet}(2739)
|                  +-{nn-applet}(2743)
|
| -nmbd(2224)
|
| -ntpd(3091)
|
| -polkitd(1197)-+-{polkitd}(1239)
|                +-{polkitd}(1240)
|
| -pulseaudio(2563)-+-{pulseaudio}(2617)
|                  +-{pulseaudio}(2623)
|
| -puppet(2373)---{puppet}(32455)
|
| -rpc.ltdn(875)
|
| -rpc.statd(954)
|
| -rpcbind(884)
|
| -rserver(1501)-+-{rserver}(1786)
|                +-{rserver}(1787)
|
| -rsyslogd(1090)-+-{rsyslogd}(1092)
|                 +-{rsyslogd}(1093)
|                 |
|                 +-{rsyslogd}(1094)
|
| -rtkit-daemon(2565)-+-{rtkit-daemon}(2566)
|                    +-{rtkit-daemon}(2567)
|
| -sd_cicero(2852)-+-sd_cicero(2853)
|                  +-{sd_cicero}(2854)
|                  +-{sd_cicero}(2855)
|
| -sd_dunny(2849)-+-{sd_dunny}(2850)
|                 +-{sd_dunny}(2851)
|
| -sd_espeak(2749)-+-{sd_espeak}(2845)
|                  +-{sd_espeak}(2846)
|                  +-{sd_espeak}(2847)
|                  +-{sd_espeak}(2848)
|
| -sd_generic(2463)-+-{sd_generic}(2464)
|                   +-{sd_generic}(2685)
```

parent and child questions...

what if parent process exits before child?

child's parent process becomes process id 1 (typically called *init*)

what if parent process never `waitpid()`s (or equivalent) for child?

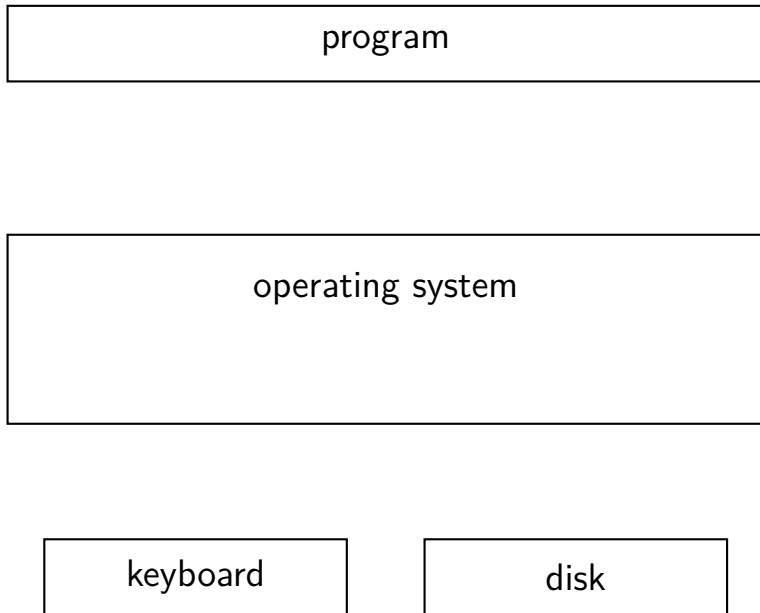
child process stays around as a “zombie”

can't reuse pid in case parent wants to use `waitpid()`

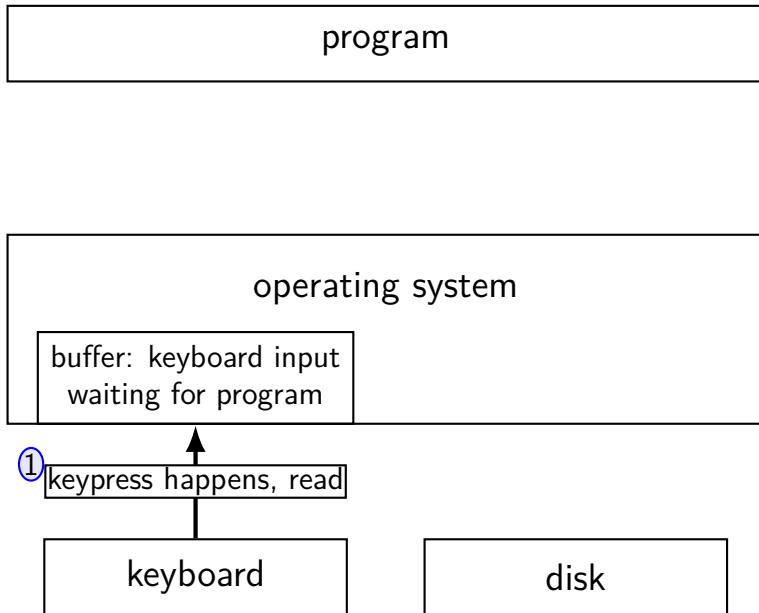
what if non-parent tries to `waitpid()` for child?

`waitpid` fails

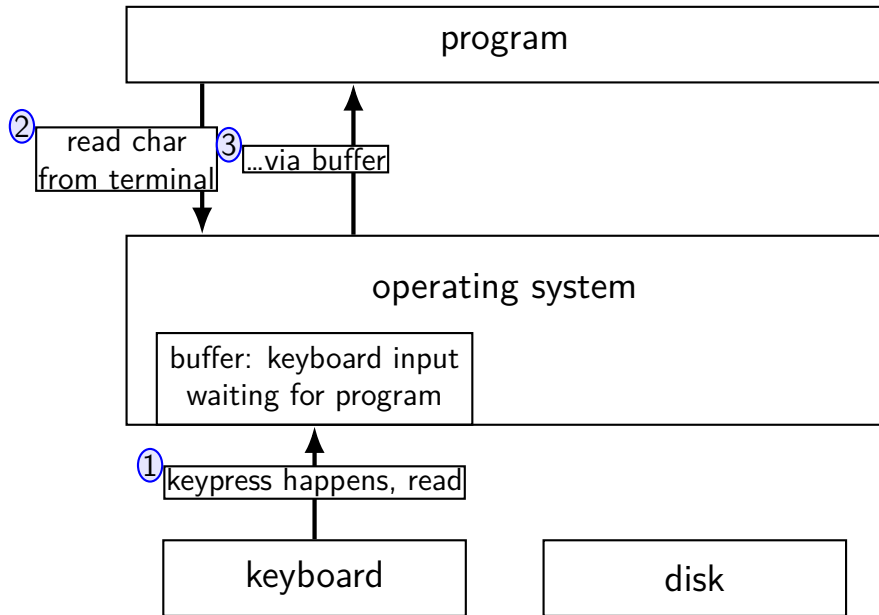
kernel buffering (reads)



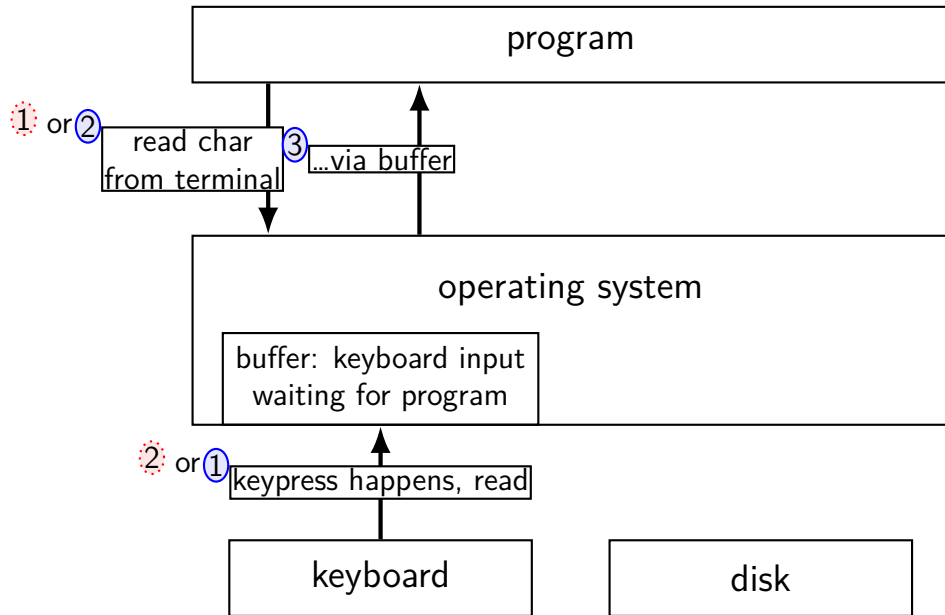
kernel buffering (reads)



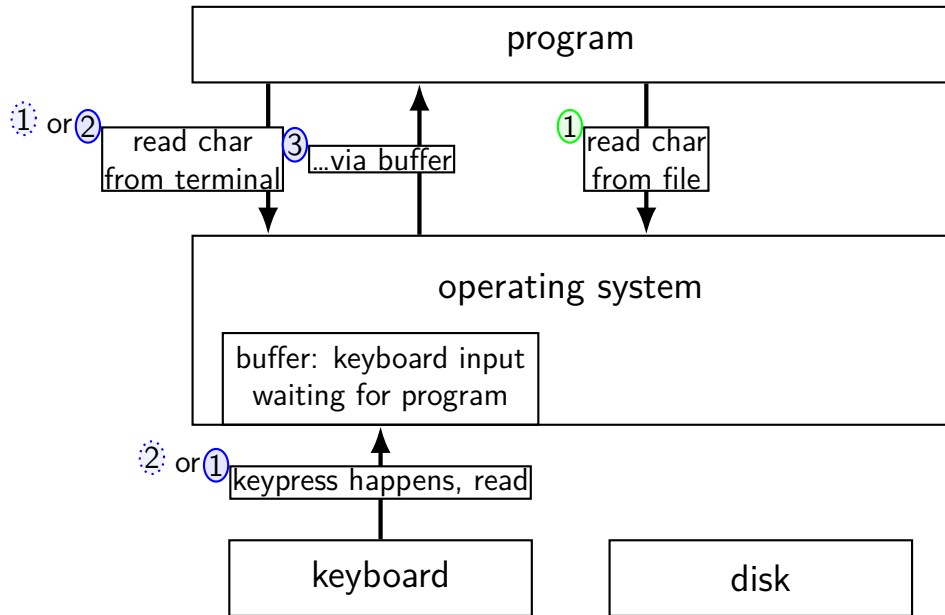
kernel buffering (reads)



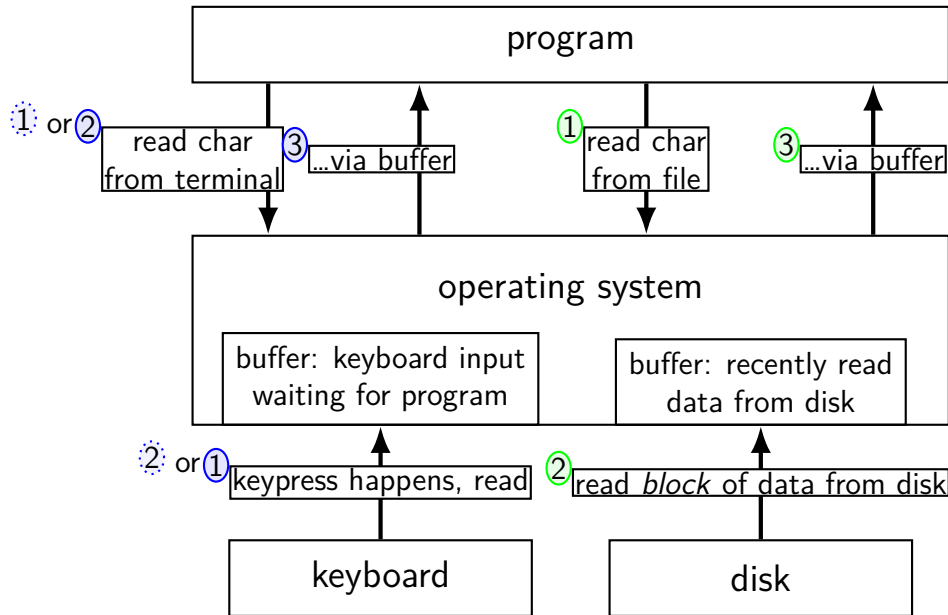
kernel buffering (reads)



kernel buffering (reads)



kernel buffering (reads)



kernel buffering (writes)



A diagram illustrating kernel buffering for writes. It consists of three main components: a 'program' box at the top, an 'operating system' box in the middle, and two output boxes at the bottom labeled 'network' and 'disk'. The 'program' box is connected to the 'operating system' box. The 'operating system' box is connected to both the 'network' and 'disk' boxes. The 'operating system' box is significantly larger than the other three boxes, indicating it acts as a central hub or buffer.

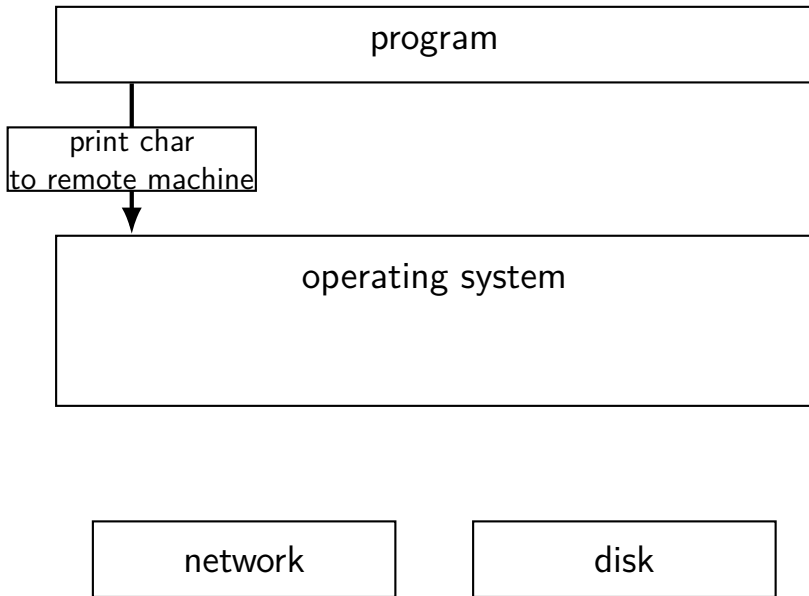
program

operating system

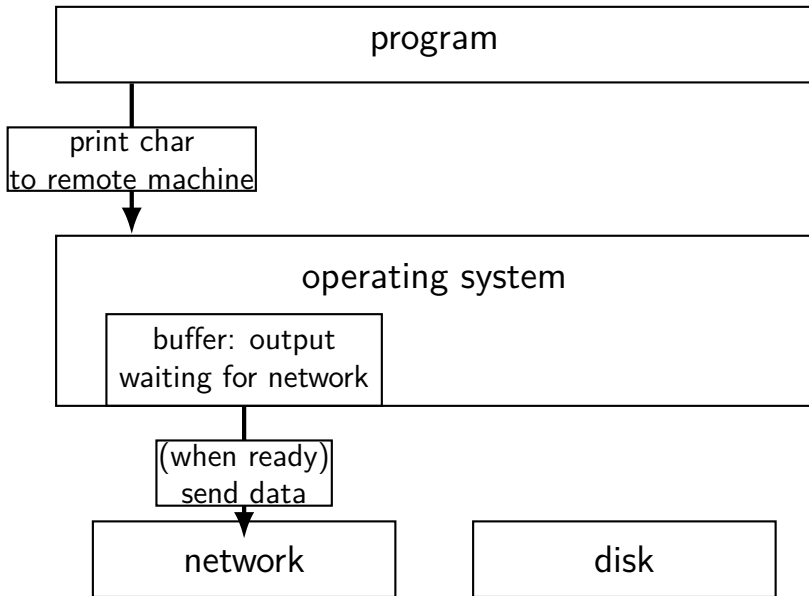
network

disk

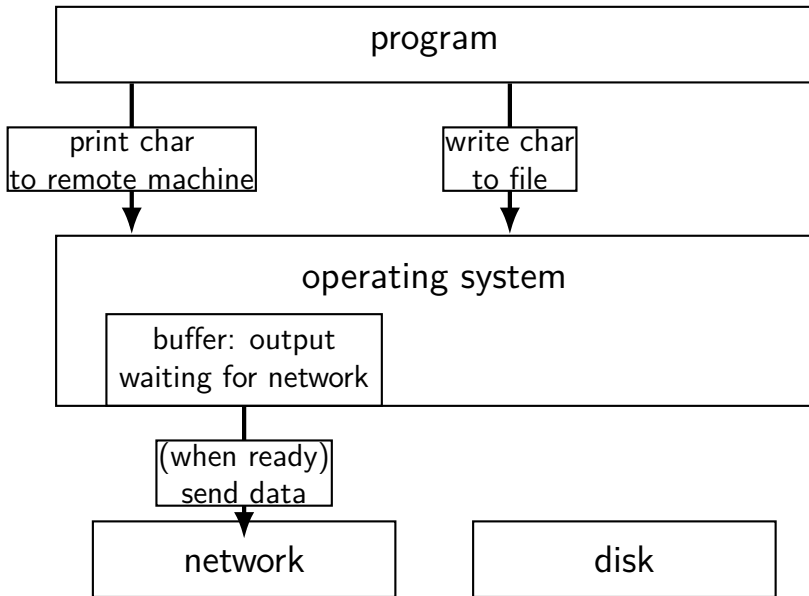
kernel buffering (writes)



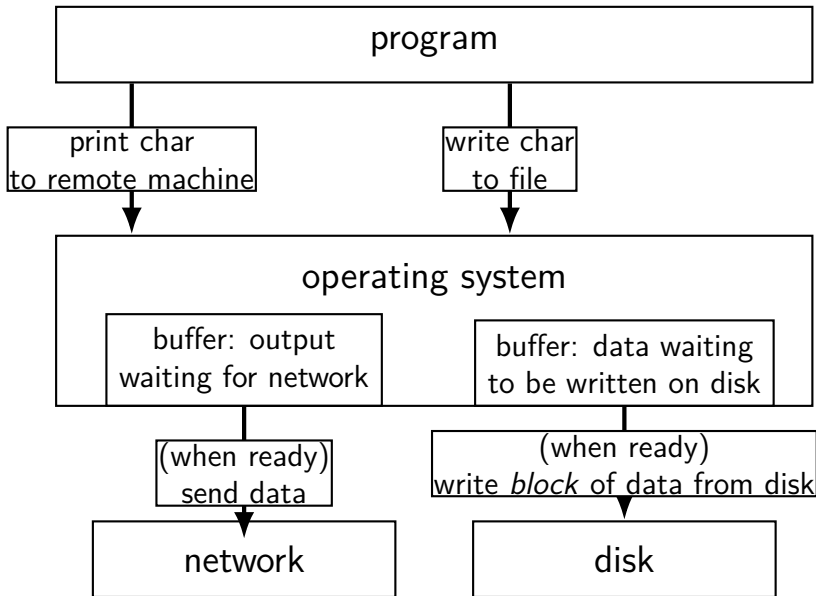
kernel buffering (writes)



kernel buffering (writes)



kernel buffering (writes)



read/write operations

read()/write(): move data into/out of buffer

possibly wait if buffer is empty (read)/full (write)

actual I/O operations — wait for device to be ready
trigger process to stop waiting if needed